Journal of Hydrology 524 (2015) 311-325

Contents lists available at ScienceDirect

Journal of Hydrology

journal homepage: www.elsevier.com/locate/jhydrol

Inductive machine learning for improved estimation of catchment-scale snow water equivalent



David Buckingham*, Christian Skalka, Josh Bongard

Department of Computer Science, University of Vermont, Burlington, VT 05405, USA

ARTICLE INFO

Article history: Received 5 December 2014 Received in revised form 18 February 2015 Accepted 21 February 2015 Available online 2 March 2015 This manuscript was handled by Konstantine P. Georgakakos, Editor-in-Chief, with the assistance of Kun Yang, Associate Editor

Keywords: Snow water equivalent Machine learning Wireless sensor network Snowpack modeling Genetic programming

SUMMARY

Infrastructure for the automatic collection of single-point measurements of snow water equivalent (*SWE*) is well-established. However, because *SWE* varies significantly over space, the estimation of *SWE* at the catchment scale based on a single-point measurement is error-prone. We propose low-cost, lightweight methods for near-real-time estimation of mean catchment-wide *SWE* using existing infrastructure, wireless sensor networks, and machine learning algorithms. Because snowpack distribution is highly nonlinear, we focus on Genetic Programming (GP), a nonlinear, white-box, inductive machine learning algorithm. Because we did not have access to near-real-time catchment-scale *SWE* data, we used available data as ground truth for machine learning in a set of experiments that are successive approximations of our goal of catchment-wide *SWE* estimation. First, we used a history of maritime snowpack data collected by manual snow courses. Second, we used distributed snow depth (*HS*) data collected automatically by wireless sensor networks. We compared the performance of GP against linear regression (LR), binary regression trees (BT), and a widely used basic method (BM) that naively assumes non-variable snowpack. In the first experiment set, GP and LR models predicted *SWE* with lower error than BM. In the second experiment set, GP had lower error than LR, but outperformed BT only when we applied a technique that specifically mitigated the possibility of over-fitting.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

There has been extensive research on techniques for measuring and modeling snow because it affects many hydrological, atmospheric, and biological processes (Tappeiner et al., 2001). The accurate estimation of snow water equivalent at the catchment scale is useful in many applications, including agricultural planning, metropolitan use, flood risk evaluation, planning of hydropower production potential, weather forecasting, and climate monitoring (Marofi et al., 2011; Schmucki et al., 2014). More than 1/6 of people globally depend on seasonal snow or glaciers for water supplies (Bales et al., 2006), and in the western United States the majority of surface water resources is derived from snowmelt (Serreze et al., 1999). However, snow has declined across much of the US over the last half-century (Pierce et al., 2008). The current severe drought in California, with record low snowpack measurements over three years, threatens water supplies throughout the state (Boxalla, 2014) and highlights the importance of snowpack research. Snow both influences climate and responds directly to climate change (Engeset et al., 2004). While climate change warrants increased snowpack monitoring, existing techniques perform poorly under extreme climatic conditions (Molotch et al., 2005; Balk and Elder, 2000), and it has been argued that the stationarity of hydrological processes can no longer be assumed (Milly et al., 2008). Furthermore, high costs of data gathering constrain the temporal and spatial granularity of estimation methods. New techniques are needed.

We propose new low-cost techniques for estimating catchment-wide snow water equivalent using machine learning algorithms, especially genetic programming. These algorithms use data gathered from existing sensor infrastructure, and possibly short-term deployments of wireless sensor networks. The manipulation of large data sets in order to gain insight into snow accumulation, melt, and runoff has been highlighted as a necessary next step in mountain hydrology (Dozier, 2011). The long-term, overarching goal of our research project is to achieve better nearreal-time (NRT), estimation of *SWE* at the catchment scale. By NRT, we mean automated reporting at fine-grained timescales, for example hourly. By better, we mean more accurate estimation without significantly increased infrastructure cost. Our strategy is to generate snow telemetry datasets using short-term, low-cost field campaigns that can be used by machine learning algorithms





HYDROLOGY

^{*} Corresponding author.

E-mail addresses: dbucking@uvm.edu (D. Buckingham), skalka@cs.uvm.edu (C. Skalka), josh.bongard@uvm.edu (J. Bongard).

to generate snowpack models. Following field campaigns and the termination of associated measurement techniques, these models can be used for NRT *SWE* estimations with no new instrumentation overhead.

The key idea behind our approach is that machine learning models are able to induce relationships between input parameters and an output value, if such exist, on the basis of the ground truth data if provided. The machine learning method we emphasize is genetic programming (GP), which generates equations relating a dependent variable to a set of independent variables.

In our case, we argue that if we obtain multiple years of "true" average *SWE* for a catchment, machine learning will be able to induce a meaningful mathematical relation between telemetry, such as proximal snow pillow reading(s), and true average *SWE*. Then, in years when true average *SWE* is not available, inputs such as snow pillow readings can be translated into average *SWE* estimates for the catchment. This approach assumes interannual continuity in snow distributions over a catchment, which has been demonstrated by previous research (Scipión et al., 2013; Tappeiner et al., 2001; Schirmer et al., 2011). Because accurate measurements of mean catchment *SWE* are generally unavailable at this time, we use snow course and wireless sensor network data as proxies for true average *SWE* to serve as ground truth for machine learning.

Thus, the ideal we aim for is a generally applicable technique for inducing models that take as input parameters existing infrastructure NRT telemetry, such as snow pillow readings, meteorological data, and date/time information, and output measurements of *SWE* at those locations. This would allow more accurate *SWE* estimation to be provided without additional cost beyond that of the initial field campaign for obtaining a ground truth dataset (Fig. 1).

Several theoretical and practical challenges exist on the way to achieving this goal. The purpose of this paper is to address them and make progress in three particular ways.

First, we explore the issue of what sort of machine learning approaches are best in this context. In general, we argue that techniques that are able to model nonlinear relationships are needed due to the known nonlinear nature of snow distribution in alpine environments (Tappeiner et al., 2001; Marofi et al., 2011). We also argue that so-called white-box tools are best, since these can provide physical insights for scientists (Schmidt et al., 2011). Furthermore, we emphasize resiliency against over-fitting, which is especially important given that the datasets available for machine learning may be relatively small.

Second, we investigate what sort of input parameters should be used by *SWE* estimation models, especially in light of practical concerns, i.e. available telemetry and datasets. In fact, availability of data is a key issue in this effort, and defines what is possible. We acknowledge the importance of terrain effects in determining snowpack distribution, influencing both accumulation and ablation patterns (Winstral et al., 2013; Fassnacht et al., 2003; Marks et al., 1999). However, because all snow sensors and courses are on flat or nearly flat ground, we did not include topographic data as explicit inputs to our models. We emphasize the flexibility of inductive machine learning, which can accommodate arbitrary new input modalities. Only those that are predictive of the dependent variable of interest will be significantly incorporated into the generated models. In this paper we focus on several potential snow telemetry and meteorological inputs in order to demonstrate the applicability of our techniques to catchment-scale *SWE* estimation, while considering the potential for future work to explore other inputs such as topographic data.

Third, we grapple with the issue of ground-truth for catchmentscale SWE and usable datasets. Constraints on our goal were imposed by the availability of snowpack data. We are not aware of catchment-wide SWE datasets with sufficiently fine time granularity to support our ideal scenario. Although datasets such as those provided by the Cold Land Processes Field Experiment (National Snow & Ice Data Center, 2014) and numerous others provide catchment-scale snowpack measurements, their time granularity is on the order of several months at least. Airborne techniques in general are cost-prohibitive for real-time reporting (Bühler et al., 2011). Although satellites are used to measure snow-covered area and albedo (Dozier and Painter, 2004), satellite retrievals of SWE are not feasible. Manual snow courses provide better temporal resolution than airborne methods (e.g. biweekly) but at low spatial resolution: snow courses measure SWE at a single location. We highlight the Snowcloud wireless sensor network, which measures HS (an effective predictor of SWE) in NRT (e.g. hourly) at multiple locations distributed over an area of interest. However, this technology is new, and available data collected by Snowcloud deployments is limited.

2. Background and contributions

Here we briefly define and summarize the machine learning methods used in this work. These techniques are described in more detail, with special emphasis on GP, in Section 4. The basic method (BM) assumes the spatial homogeneity of SWE. It naively estimates mean catchment-wide SWE to be the same as the single-point SWE measurement taken at a snow pillow. Linear regression (LR) fits a least-squares linear model to training data (Hastie et al., 2009). The prediction is a weighted linear combination of the input variables. Binary regression trees (BT) are nonlinear models which are generated using training data (Hastie et al., 2009). A BT model partitions a set of predictions according to the input variables such that a given set of input values results in a specific prediction. Genetic Programming (GP) is a symbolic regression algorithm that uses training data to iteratively improve a population of nonlinear models through a combination of stochastic variation and performance-based selection (Koza, 1992).



Fig. 1. First, the Snowcloud WSN is deployed in an area near a snow pillow. Next, data generated by Snowcloud, by the pillow, and potentially other sources, are used by machine learning to generate a model of snowpack distribution. Finally, after Snowcloud has been removed, the model is used to estimate snow levels in the area where Snowcloud had been deployed.

In our ideal situation we would use a large set of accurate measurements of mean catchment SWE as ground truth to train and evaluate models that predict mean catchment SWE in NRT. However, the only SWE measurements available at this spatial scale are generated by airborne techniques with time resolutions that are insufficient for machine learning (e.g. twice per year). Because machine learning needs a large number of samples for model training and because we want to predict SWE in near-realtime, we required much more frequent measurements. We therefore developed a series of experiments using available snowpack data in lieu of NRT catchment-scale SWE measurements to explore successive approximations of our ideal scenario. Approximations of average catchment SWE, obtained via snow courses and distributed ground-based sensor readings, serve as ground truth for machine learning in our experiments. Implicit in our work is the importance of new methods for obtaining NRT catchment-scale SWE ground-truthing via low-cost distributed sensor networks. As data from NASA's Airborne Snow Observatory (NASA Airborne Snow Observatory, 2015) become available for a range of years, they will provide an ideal data set for our approach.

First, we used snow course measurements, which involve the manual collection of SWE and/or HS at a single location, as a proxy for catchment-wide SWE. Although snow courses do not directly measure snowpack distribution at the catchment scale, they are likely to provide measurements that are *closer* to mean catchment SWE than snow pillows measurements are. Snow courses take multiple measurements over approximately 200 m, so they involve a much larger sample size than the single-point measurements of snow pillows. Furthermore, pillow under-measurement or overmeasurement errors may occur when the base of the snow cover is at melting temperature (Johnson and Marks, 2004). Thus, we used snow course data as a first approximation of mean catchment SWE to provide ground-truth data for machine learning. We generated models that use readily available information such as meteorological telemetry and snow pillow measurements as input variables. This approach, which is explored in Experiment Set I. would allow for shorter or less frequent snow courses or for their discontinuation and, because it uses previously collected data. incurs no data gathering costs.

Second, we used *HS* data collected by the Snowcloud (Skalka and Frolik, 2014) wireless sensor network (WSN) at sites in Norway and California, each for only one snow season, as a proxy for catchment-wide *SWE* data. Snowcloud is a WSN-based data gathering system for snow hydrology, notable for its low-cost and ease of deployment, developed and operated by the University of Vermont. A network of light-weight sensor towers (nodes) is deployed over an area of interest for a short-term field campaign to collect spatially distributed measurements of relevant meteorological processes (Fig. 3). In addition to *HS*, Snowcloud measures air temperature, soil temperature, and solar radiation. Mesh wireless communication allows data from the entire network to be collected wirelessly by communication with a single node.

We used measurements collected from Snowcloud over the course of a single snow season to generate ground-truth estimates for model-training. Note that it could be desirable to collect data over multiple seasons as models trained on multi-year data may be more robust against internal-annual variations in snowpack distribution. Once a model has been obtained, the WSN may be recovered for re-deployment at another site. Unlike pillows and snow courses, Snowcloud collects NRT data from multiple locations, potentially capturing more of the variability of snowpack distribution than is possible with single-location measurements. Thus, we use Snowcloud data as a second approximation of catchment mean *SWE* to provide ground-truth data for machine learning. This technique is explored in Experiment Set II.

Recent research by Kerkez et al. (2012) and Welch et al. (2013) has developed new sensor placement strategies for monitoring snow. Although these methods were not employed in the experiments discussed in this paper, they should be considered in future applications of our techniques.

2.1. Suitability of machine learning

Snow pillows are large, expensive, permanent installations that measure *SWE* at a single location. The infrastructure for the automatic collection of *single-point SWE* is well established. For example, there are 830 Snowpack Telemetry (SNOTEL) sites in the United States (Surveyor, 2014) and another 124 snow pillows operated by the California Department of Water Resources. However, the extrapolation from single-point measurements to surrounding areas is error prone. The spatial distribution of alpine snow cover is highly variable (Balk and Elder, 2000; Elder et al., 1991; Jost et al., 2007), due to a variety of environmental forcing effects, such as topography (Anderton et al., 2004), canopy cover (Moeser, 2010), and wind and solar exposure (Moeser, 2010; Moeser et al., 2011).

Meromy et al. (2013) studied 15 snow stations across the western United States and found that snow station biases were frequently greater than 10% of the surrounding mean observed snow depth. The flat-field areas where snow pillows are commonly located are usually not typical of more complex nearby terrain, causing the majority of such stations to overestimate snow depth in their vicinity (Grünewald et al., 2013). Molotch and Bales (2005) studied the areas surrounding six SNOTEL stations in the Rio Grande headwaters. They found that only a small fraction of grid elements were representative of mean grid SWE during accumulation, and that no elements were representative of mean grid SWE during both accumulation and ablation. SNOTEL stations in the Rio Grande headwaters preferentially represent densely forested areas and experience snow cover persistence that is 14% greater than the mean persistence of the watershed (Molotch and Bales, 2006). Rittger (2012) found that errors based on statistical relationships between point measurements of snow and streamflow in the Sierra Nevada can reach 25-70% in one out of five years.

The relative importance of separate processes which govern snow distribution varies over the course of a snow season. Elder et al. (1991) summarize the various processes and explain how their influence changes over time. During the winter, accumulation and redistribution processes dominate. Precipitation is determined by regional climate and latitude as well as by local orographic effects, and redistribution by wind, avalanches, and sloughs are the primary causes of spatial heterogeneity. In the spring, however, snow distribution is controlled mainly by ablation. Of the many energy sources, solar and longwave radiation dominate. This energy decreases water in a basin through sublimation and when runoff leaves the basin. It also redistributes SWE, affecting spatial variability. These dynamics highlight the need for NRT modeling of snowpack, as the forcing effects that establish snow distribution vary drastically over the course of a snow season.

However, the significant *consistency* of snowpack *between* years encourages investment into the development of reusable statistical models. Strong inter-annual consistency in the spatial distribution of snow (Scipión et al., 2013), in SCA (Tappeiner et al., 2001), and in the snow depth patterns of maximum accumulation (Schirmer et al., 2011), have been observed in the Swiss and Italian Alps. In the western United States, consistent wind directions can produce stable snow accumulation patterns from year-to-year (Winstral and Marks, 2014). These findings suggest a strong link between accumulation patterns and geophysical terrain and indicate that site-specific snow distribution models may be able to accurately characterize snowpack distribution over multiple years. Nevertheless, long-term changes in the patterns of snow distribution may be caused by factors such as changes in vegetation or climate change. Therefore, it may occasionally be necessary to rerun GP and generate a new model. Techniques such as retroactive *SWE* calculation (Rittger et al., 2011) could be used to detect when previous models begin to perform poorly, indicating that secular variability in the dynamics of snow distribution warrants the development of a new model.

It may be desirable to produce non-site-specific models. Trained at catchments where ground truth data is available, and making use of predictor variables that vary between catchments, such as topography, such models could then be applied to catchments where no independent measurement of mean catchment SWE exists. However, we did not incorporate topography because the snow pillows are all on flat or nearly flat ground. Our work focuses on site-specific models and use model inputs that vary over time at a given catchment.

2.2. Why GP?

It has been demonstrated that the relationships between snow distribution and the topographic and meteorological forcing effects include nonlinearities (Tappeiner et al., 2001), and the spatial distribution of *SWE* is nonlinear because it is influenced simultaneously by numerous processes including accumulation, ablation, and snow drifting (Marofi et al., 2011). GP can produce both linear and nonlinear models. If the data used to train GP contain only linear relationships, the resulting models will be linear, and the performance of GP will be similar to that of LR.

White-box models, such as those produced by GP, can be interpreted by human analysis, potentially yielding new information about the modeled data (Schmidt et al., 2011). Some nonlinear regressors, such as artificial neural networks, produce models that are difficult or impossible to interpret. GP trees, however, can be expressed as mathematical equations (Fig. 2). It is possible that by examining these equations domain experts could gain novel insight into the processes governing snow distribution.

Unlike regression techniques that constrain the form of the regressor, GP can combine operators, variables, and constants into arbitrary arrangements. GP does not require any assumptions about the form that a model should take: it is left open to inductive search. By generating models that use predictor variables in unexpected ways, GP may help discover previously unknown relationships among variables.



Fig. 3. Snowcloud WSN sensor tower. A complete sensor stand with solarrecharged battery power, wireless mesh communication, and multiple sensor modalities. October 2011, Mammoth Lake, CA.

Finally, as we will discuss further, GP may be augmented with multi-objective optimization, which constrains GP to produce parsimonious models. This mitigates against over-fitting, a significant concern in the case that relatively small datasets are available for machine learning.

While many regression techniques possess one or more of these desirable qualities, GP possesses all of them, making it an ideal candidate for snowpack modeling.

2.3. The primacy of snow depth

While *SWE* is a product of *HS* and density (ρ), it has been shown that *HS* is the essential determining metric for *SWE* estimation. Models have been developed to derive ρ estimates from *HS* measurements (Logan, 1973; Sturm et al., 2010), and measurements of *HS* are highly predictive of *SWE* (Adams, 1976). Analysis of the spatial variability of *HS* and ρ has revealed that the variability of *HS* is significantly greater than that of ρ (López-Moreno et al.,



314

Fig. 2. These example GP trees were manually selected from the final populations of GP runs conducted for Experiment Set II. The leftmost tree represents a simple linear model. The middle tree is a nonlinear model. The rightmost tree is a more complex nonlinear model.



Generation 0

 $y = (log(x) + 8.293)^{-2}$ y = sin(x) + 0.388 $y = (-x - 0.319)^{x}$ $y = 1.303 * x^{(x^{1.07})}$

Generation 1

y = sin(x) + 0.388 y = sin(x - 0.026) + 0.388 $y = 1.303 * x^{(x^{1.07})}$ $y = 0.912 * x^{(x^{1.07})}$: Generation n (.0.501)

 $y = \cos(x * 1.309) - (x^{0.501})$ y = ((x - 0.026) * 1.204) + 0.388 $y = (0.912 * x^{(x^{1.81})}) - 0.441$ $y = (7.337 * (x^{1.81})) - 8.139$

Fig. 4. Genetic programming algorithm. The figure on the left demonstrates the iterative process through which GP modifies a population of solutions. On the right, a population of four models evolves as each iteration of the GP cycle produces a new generation.

2012). Variation of *SWE* is therefore overwhelmingly a product of *HS* variation (Moeser et al., 2011; Molotch et al., 2005; Sturm et al., 2010; Elder et al., 1991, 1998). The effect of ρ variation on *SWE* is small by comparison, and estimates of areal *SWE* derived from one or several *SWE* measurements can be greatly improved by incorporating a larger number of *HS* measurements (Elder et al., 1998; Moeser et al., 2011), which are much less labor intensive than manual *SWE* measurements (Sturm et al., 2010). Snowcloud, which provides ground-truth data Experiment Set II, measures *HS*. Therefore, as has been done elsewhere, we use *HS* as a "surrogate for *SWE*" (Winstral et al., 2002).

2.4. Related work

Moeser et al. (2011) explored three models for estimating *SWE* in the area around a meteorological station using ground based measurements. The first model used meteorological data such as air temperature and solar radiation, tree canopy cover measurements, and *HS* measurements collected by the Snowcloud WSN, as well as a single-point *SWE* measurement. The second model used multiple *HS* measurements and single-point *SWE* measurements, but no meteorological or tree canopy data. The third model used meteorological and tree canopy data, along with multiple *HS* measurements, but no single-point *SWE* measurement. It was found that increasing the number of *HS* measurements can improve areal *SWE* measurements because *HS* varies more than snow density. While this work used linear modeling; our work expands upon it by developing nonlinear models.

Marofi et al. (2011) compared three methods for modeling *SWE*: multivariate nonlinear regression (MNLR), artificial neural networks (ANN), and a neural network-genetic algorithm (NNGA), where genetic algorithms were used to parameterize ANNs and the learning process. ANN performed better than MNLR, suggesting that computational intelligence approaches may outperform MNLR for modeling *SWE*. NNGA performed better than ANN, suggesting that evolution-inspired genetic algorithms can be used to develop effective models of *SWE*. Tabari et al. (2010) estimated *HS* and *SWE* using multiple methods and also found that NNGA provided the best results. Unlike neural networks, GP produces white box models.

Tappeiner et al. (2001) compared the performance of LR-based and ANN-based snowpack models, which used topographic and meteorological data to estimate *SWE*. The authors compared the results of LR with ANN to estimate the degree of necessary nonlinearity in *SWE* modeling. The ANN performed significantly better than LR, demonstrating nonlinearity in the relationships between topographic and meteorological variables and *SWE*.

Several studies have used binary regression trees to model snowpack. Winstral et al. (2002) derived terrain-based parameters from digital elevation models (DEM) which were used as input variables to binary regression trees. One parameter was based on maximum upwind slopes relative to seasonally averaged winds. Another measured upwind breaks in slope from a given location. Binary tree models based on these terrain-based parameters as well as elevation, solar radiation, and slope performed better than models based only on elevation, solar radiation, and slope, Elder et al. (1998) modeled the distribution of SWE by merging remotely sensed snow-covered area data with binary tree models applied to field measurements of HS and SWE. Balk and Elder (2000) combined binary regression trees with kriging of manual snow survey measurements and snow-covered area determined by aerial photographs, to estimate SWE. Anderton et al. (2004) used binary regression trees to relate HS and disappearance date to terrain indices. They found that the topographic effects on snow redistribution by wind primarily determined SWE distribution at the start of the melt season which, more than melt rates, determined the patterns of snow disappearance. Molotch et al. (2005) compared binary regression tree models using various sources of DEMs and found that using DEMs from different sources leads to significant differences in modeled snowpack distribution. The most significant differences were on ridge-tops, where the elevation values differed across DEMs.

In Experiment Set II we compare the performance of BT to GP. Unlike this previous work which used binary regression trees to produce spatially distributed models of snowpack, our models predict a single value: mean *HS* measured by a wireless sensor network.

Marks et al. (1999) also developed spatially distributed models. They used topographic data to determine estimates of radiation, temperature, humidity, wind, and precipitation for use in a coupled energy and mass-balance model called ISNOBAL.

Recent research has made significant advances in simulating the effects of wind on snow distribution. Winstral et al. (2009) developed a simplified wind model that uses upwind topography to accurately predict wind speeds. Winstral et al. (2013) developed

Table 1	
CDEC snow course site descriptions.	

ID	EL (m)	Name	Asp.	Exposure
CAP	2438	Caples Lake	SW	open meadow, low brush
\mathcal{GRZ}	2103	Grizzly Ridge	Ν	meadow in scattered timber
KTL	2225	Kettle Rock	S	sloping, open meadow
MSH	2408	Mount Shasta	SE	grassy and rocky meadow
\mathcal{NTH}	2835	North Lake	SE	grassy meadow
SPD	1585	Lake Spaulding	level	grassy meadow
\mathcal{HIG}	1838	Highland Lakes	NW	medium sized meadow in dense timber
HYS	2012	Huysink	W	open meadow on one leg, opening in timber on second leg

a snow distribution algorithm that uses terrain structure, vegetation, wind, and precipitation data to simulate wind-affected snow accumulation. It accurately predicted disparate snow distribution caused by inhomogeneous precipitation and redistribution by wind. Winstral and Marks (2014) analyzed the effects of wind on snow distribution. They found that high wind speeds increased snow depth variability, that forested sites decreased variability by moderating wind effects, and that consistent wind directions produced accumulation patterns that were stable between years.

Sturm et al. (2010) used *HS*, day of the year, and climate classes, such as Alpine, Maritime, and Tundra, to estimate snowpack density. Estimated snowpack density was used to convert *HS* measurements into *SWE* estimates.

Guan et al. (2010) found that atmospheric rivers (ARs), are associated with intense storms that contribute a large percentage of snow during most years. Because AR storms are relatively warm, the participation of AR participation into snowfall versus rainfall is sensitive to minor variation in surface air temperature.

Rittger et al. (2011) combined satellite-based measurements of snow-covered area with energy balance calculations to retroactively calculate distributed SWE at the date of maximum accumulation, using the "reconstruction" technique originally developed by Martinec and Rango (1981). This calculation was then used to evaluate the accuracy of two real-time models. They found that at elevations below 1500 m, the real-time models overestimated *SWE* because of early season melt, and at elevations above 3000 m, the real-time models underestimated *SWE* because they do not sample these higher elevations. It is possible that this technique could be used to evaluate the effectiveness of the inductive learning methods that we describe in this work.

3. Training data and model inputs

Inductive machine learning requires substantial datasets for developing and evaluating models, and we acquired extensive hydrological and meteorological data for use in our experiments. We focused on two types of available datasets that are approximations of mean catchment SWE. First, we consider a record of CDEC snow courses from the Sierra Nevada. We observe that CDEC snow courses are intended to provide an estimation of SWE at a particular elevation (USDA, 2014), though in fact they are linear transects of SWE samples. Second, we consider a record of Snowcloud sensor network readings from Norway and California. Snowcloud provides distributed coverage of snow depth readings for the deployment area, as well as fine time granularity, and can support better estimations of mean catchment *SWE* than periodic snow courses.

3.1. Experiment Set I data

Experiment Set I used data collected from eight sites across California. There were three main types of data: *SWE* from manual snow courses, *SWE* measurements from snow pillows, and air temperature data.

The California Data Exchange Center (CDEC) provided an extensive database of snow data. The snow courses that we used, which are described in Table 1, were performed monthly, were about 200 meters long, and consisted of 10 measurements, the mean of which was recorded. CDEC also maintains single-point *SWE* measurement data from snow pillows at sites throughout California. Of the 404 snow course sites, 59 are co-located with snow pillows.

The National Climate Data Center (NCDC) maintains meteorological data, such as air temperature, wind speed, and solar radiation measurements, collected at weather stations across the United States. We used data from the four NCDC stations which are located within 30 km of CDEC snow courses. We arbitrarily chose a 30 km cutoff because we suspected that meteorological activity within that distance might be predictive of measurements at the snow course. The models generated by machine learning will not make significant use of input data that is not predictive.

Significant gaps exist in the NCDC database, and of the various sensor modalities, air temperature data is the most complete. Using more meteorological inputs and necessarily fewer data samples, we had previously been unable to generate effective models of *SWE*. For Experiment Set I, therefore, air temperature was the only meteorological input. Air temperature is known to be a highly effective predictor of melt rate because it is correlated with longwave atmospheric radiation, the most important energy source for snowmelt (Ohmura, 2001). Air temperature is made accessible to the models by three variables: *minTemp7*, *maxTemp7*, and *meanTemp7*, which aggregate daily values over the seven days inclusively preceding the day for which *SWE* is estimated.

We used the temporal and spatial intersection of available data from these three sources (CDEC snow courses, CDEC snow pillows, NCDC air temperature data) to construct eight datasets, based on eight snow course sites. These snow courses were selected because they are coincident with either snow pillow data, NCDC air temperature data, or both, over a range of time that includes a large number of samples points (greater than 100 except for one site). The constructed datasets are summarized in Table 2.

3.2. Experiment Set II data

Experiment Set II used *HS* data collected by four Snowcloud sensor nodes in Sulitjelma, Norway between January and April, 2013. Each node sampled *HS* every six hours. We averaged *HS*

able 2					
xperiment Set I	data	summary	by	CDEC	site.

ID	Pillow	NCDC base	Dist (Mi)	Samples	Years
\mathcal{CAP}	YES	N/A	N/A	177	1970-2011
GRZ	YES	N/A	N/A	207	1970-2011
\mathcal{KTL}	YES	N/A	N/A	159	1979-2011
\mathcal{MSH}	NO	Mount Shasta	5.98	137	1973-2011
\mathcal{NTH}	NO	Bishop Airport	18.27	147	1973-2011
SPD	NO	Blue Canyon Nyack	4.56	174	1977-2011
\mathcal{HIG}	YES	Mount Shasta	18.31	75	1980-2012
HYS	YES	Blue Canyon Nyack	9.79	111	1984-2011

Table 3Snowcloud deployment coordinates.

Sulitjelma, Norway			Sagehen,	CA	
Tower	Lat.	Long.	Tower	Lat.	Long.
1	67.0981	16.0488	1	39.43161	-120.23975
2	67.0983	16.0497	2	39.43155	-120.23936
3	67.0983	16.0482	3	39.43140	-120.23976
4	67.0987	16.0487	4	39.43173	-120.23882
			5	39.43173	-120.23864
			6	39.43204	-120.23872

measurements from the four nodes (Table 3) and then over each day to produce 93 estimates of mean catchment *HS*. These values served as ground-truth *HS* for experiments at Sulitjelma.

Approximately 16 km away from the Sulitjelma Snowcloud deployment site is Storstilla nedanför Balvatn in Nordland County, station number 164.12.0 (Balvatn). The Balvatn station records both *HS* and *SWE*. Daily *HS* measurements collected at Balvatn compose the *HS* input variable to models developed for Sulitjelma in Experiment Set II.

Six Snowcloud wireless sensor network sensor nodes were deployed within the Sagehen Creek Field Station, near Truckee, California, from January to May, 2010. Each node reported daily *HS* measurements, which we averaged to generated 99 estimates of mean catchment *SWE*. These values served as ground-truth *HS* for experiments at Sagehen. Note that the same WSN data were used by Moeser (2010).

In order to assess the significance of the *source* of single-point *HS* input variables, we developed models for estimating mean *HS* at the Sagehen Snowcloud deployment using inputs from two different CDEC sites, *Independence Camp* (*TDC*) and *Huysink* (*HYS*). *TDC* is approximately 5.5 km away from the Snowcloud deployment and, like Sagehen, is on the Eastern side of the Sierra crest. *HYS* is approximately 30 km away, on the Western side of the crest.

3.3. Time of year

Because the dynamics underlying snowpack distribution vary over the course of a snow season, for example between periods dominated by deposition and periods dominated by ablation, we introduce time of year (*TOY*) as an independent variable for both experiment sets. This allows models to distinguish parts of the snow season. Time of year is an integer value expressing the number of days since January 1.

3.4. Preparation of datasets

We define a dataset, *D*, for each experiment (each row of Table 6 and each location in each row of Table 5). Elements of a dataset *D*

(a) Random division: dataset is randomly divided into three subsets of equal size.

(c) Three bins: dataset is divided into three temporally contiguous bins, which are each divided into three subsets.

(e) Three bin case illustrating random offset. take the form of a 3-tuple, $\langle T, \theta, \vec{p} \rangle$, where *T*, time, specifies a calendar date, θ is an estimate of the true value of the independent variable, and \vec{p} is a vector of predictor variables. Although *T* is used to generate predictor variables such as *TOY* and air temperature statistics, it is not itself a predictor variable and is therefore not included in \vec{p} . *T* is unique in *D* so that no two data samples in *D* have the same *T*:

$$\forall \langle T_1, \theta_1, \vec{p}_1 \rangle, \langle T_1, \theta_2, \vec{p}_2 \rangle \in D \qquad \theta_1 = \theta_2 \qquad \text{and} \qquad \vec{p}_1 = \vec{p}_2 \qquad (1)$$

In Experiment Set I, θ is an approximation of mean catchment *SWE* derived by manual snow course. In Experiment Set II, θ is an approximation of mean catchment *HS* derived from Snowcloud WSN measurements.

Depending on the experiment, \vec{p} includes some combination of *HS* measured at a snow pillow, *SWE* measured at a snow pillow, *TOY* (an integer value derived from *T*), and air temperature, (which is composed of three variables: *minTemp7*, *maxTemp7*, and *meanTemp7*). The *Model inputs* columns of Table 5 and Table 6 specify the contents of \vec{p} for each experiment.

In order that a model developed from *D* may be evaluated on new, unseen data, *D* is divided into training, ϱ , and testing, τ , subsets. The training set is twice as large as the testing set. However, GP and BT require that ϱ be further divided into grow, *g*, and selection, *s*, subsets:

$$\varrho = g \cup s$$
 and $g \cap s = \emptyset$ and $|g| = |s|$ (2)

In all experiments, *D* is first divided into *g*, *s*, and τ :

$$D = g \cup s \cup \tau$$
 and $g \cap s \cap \tau = \emptyset$ and $|g| = |s| = |\tau|$ (3)

For BM and LR, g and s are simply combined into ϱ and used as training data. As discussed in more detail in Section 4, in the case of GP and BT g is used to generate a set of models and s is used to determine which one should be kept and evaluated on τ . In any case, ϱ is used to obtain a single model, which is then exposed to τ to evaluate its ability to predict unseen data.

We explored several methods for dividing *D* into *g*,*s*, and τ . In Experiment Set I and in the first part of Experiment Set II (Experiment Set II: *Random Division*), the chronologically ordered *D* is randomly shuffled and then divided into thirds, as illustrated by Fig. 5a. This method has the effect that a large portion of the training data is likely to be temporally proximal to testing data.

As discussed further in Section 5, we found in Experiment Set II that the temporal proximity between ϱ and τ caused machine learning to map *TOY* values to estimates of *HS*. The models memorized the data rather than capturing the relationships among the data. We therefore conducted Experiment Set II: 4 *Bins*. Instead of shuffling *D*, we maintained its ordering and divided it into four



(b) Four bins: dataset is divided into four temporally contiguous bins, which are each divided into three subsets.

۰.			

(d) Two bins: dataset is divided into two temporally contiguous bins, which are each divided into three subsets.

chronologically contiguous bins. Each bin is then subdivided into three chronologically contiguous subsets which are assigned to g, s, and τ . This method is illustrated by Fig. 5b. We also conducted Experiment Set II: 3 Bins and Experiment Set II: 2 Bins, as illustrated in Fig. 5c and d. As we move from Experiment Set II: Random Division to Experiment Set II: 2 Bins, the division of D transitions from finer to coarser temporal granularity. As this granularity becomes coarser, it becomes more difficult for machine learning to use TOY to simply memorize data. However, it also becomes more difficult for models to capture the variation of the dynamics of snowpack distribution over the course of a snow season.

In order to introduce stochasticity into the division *D* and thus allow the repetition of experiments to produce a distributed sample of results, a randomly generated offset shifts the starting point of the division. Fig. 5e illustrates the effect of this offset in the case of three bins.

4. Calculation

In this section we first describe how we compared the performance of different snowpack modeling techniques. We then describe the various modeling techniques that we used, with special emphasis on GP.

4.1. Comparing estimation methods

In order to compare the performance of two machine learning techniques, *M* and *M'*, on a dataset *D*, *D* is divided into complementary subsets ϱ and τ . Methods *M* and *M'* are applied to ϱ to produce estimators $\hat{\theta}$ and $\hat{\theta}'$. This process may be deterministic or nondeterministic. In Experiment Set I and Experiment Set II: *Random Division*, nondeterminism is introduced by the random division of *D*. GP introduces further nondeterminism by the stochasticity of the GP algorithm. The BT algorithm is deterministic when a single input variable is used, but nondeterministic when applied to multiple input variables. Estimators $\hat{\theta}$ and $\hat{\theta}'$ are applied to τ to determine the mean absolute errors of the estimators MAE($\hat{\theta}$) and MAE($\hat{\theta}'$), as we will discuss in Section 4.2.

This process of randomly dividing *D* and applying *M* and *M*^{*t*} to obtain MAE($\hat{\theta}$) and MAE($\hat{\theta}'$) is repeated 30 times, resulting in vectors of estimator errors \vec{e}_M and $\vec{e}_{M'}$ each with cardinality 30. We consider \vec{e}_M and $\vec{e}_{M'}$ to be statistical samples of errors drawn from the population of errors that method *M* and *M'* could produce given *D*. We chose to collect 30 samples because a sample size of at least 30 allows the Central Limit Theorem to be safely applied without assuming a normal population distribution, permitting the application of the one-sample *t*-test to calculate confidence intervals and the paired two-sample *t* test to test hypotheses.

The means of \vec{e}_M and $\vec{e}_{M'}$ are unbiased estimates of the true population means μ_M and μ'_M . To find out if M' outperforms M on dataset D we pose the hypotheses:

$$\begin{array}{l} H_0: \mu'_M = \mu_M \qquad (Null hypothesis) \\ H_a: \mu'_M < \mu_M \qquad (alternative hypothesis) \end{array}$$

and apply the Student's *t*-test for paired samples to \vec{e}_M and $\vec{e}_{M'}$. If the Null hypothesis is rejected, we say that method M' produces lower error (performs better) on dataset *D* than does *M*. We report the *p*-value, the probability that we have performed a Type I error by rejecting a true Null hypothesis.

4.2. Evaluating estimator error

Recall that an element *d* of dataset *D* takes the form $\langle T, \theta, \vec{p} \rangle$ and that *D* has been divided into ϱ and τ . An estimation method *M* is

applied to $\varrho \subset D$ to generate an estimator $\hat{\theta}$, which is a function from predictor variables \vec{p} to dependent variable *y*, an estimate of θ .

$$\hat{\theta}: \vec{p} \to y \qquad y \approx \theta$$

The error of $\hat{\theta}$ on an input vector is the difference between the estimate it produces and ground truth.

$$\mathbf{E}_{\hat{\theta}}(\vec{p}) = \hat{\theta}(\vec{p}) - \theta \tag{4}$$

The error is calculated on each sample in τ to determine the mean absolute error of the estimator:

$$\mathsf{MAE}(\hat{\theta}) = \frac{\sum_{i=1}^{k} |\mathbf{E}_{\hat{\theta}}(\vec{p}_i)|}{k} \tag{5}$$

where

4.3. Basic method

The basic method (BM) assumes that SWE as measured at a snow pillow is representative of catchment-wide SWE. It naively estimates ground truth (snow course-derived) SWE to be the same as the independent variable (snow pillow-derived) SWE measurement. Error in the predictive power of BM expresses the difference between snow pillow measurements and snow course SWE measurements. If *x* represent SWE measured at the snow pillow, then

$$x \in \vec{p}$$
 and $\hat{\theta}(\vec{p}) = x$ (6)

Unlike the more sophisticated machine learning techniques, BM does not make use of training data to generate a model.

4.4. Linear regression

Linear regression (LR) fits a least-squares linear model to training data which is then evaluated on test data (Hastie et al., 2009). LR expresses the linear relationships between independent and dependent variables. We used the *gsl_multifit_linear* function from the GNU Scientific Library (GSL, 2014) to perform LR. We include LR in order to gain insight into the data we are using. LR will perform less well than nonlinear techniques only if the modeled data contain nonlinear relationships.

4.5. Genetic programming

GP is an evolutionary algorithm, inspired by biological evolution, that iteratively evolves populations of parse trees to perform symbolic regression (Koza, 1992) (see Fig. 4). In this work, the trees are snowpack models, estimator functions, that use available independent variables to estimate mean *SWE* (Experiment Set I) or *HS* (Experiment Set II) at the catchment scale. Tree terminals are input variables and constants, while internal nodes are arithmetic operators. The operators we used are listed in Table 4.

Tal	ole 4	
GP	parameters	

Ĩ	
Parameter	Value
population size number of generations max tree size mutation operators binary operators unary operators terminals	1000 (Experiment Set I), 2000 (Set II) 3000 (Experiment Set I), 10,000 (Set II) 30 crossover (60%), mutation (40%) addition, subtraction, mult., division, power log, exponential, sine, cosine, independent variables, constants values

We used the lil-gp Genetic Programming System (System, 2013), an open source implementation of GP, in order that we might make any needed modifications. We modified lil-gp to implement multi-objective Pareto optimization.

GP begins by generating a starting population of randomly constructed trees. Each tree in the population is evaluated on training data to determine its fitness, defined as the inverse of mean error. Trees are selected according to their size and fitness to produce the population for the next generation. Genetic operators make stochastic modifications to the new trees, randomly perturbing their fitness values. The genetic operators we used were mutation and crossover. Mutation, which is applied to 40% of new trees, selects a subtree at random and replaces it with new, randomly generated subtree. In crossover, which is applied instead of mutation 60% of the time, two parent trees exchange subtrees, resulting in two novel offspring. Crossover allows recombination of subtrees from existing models while mutation introduces new subtrees to the population, maintaining genetic diversity. Because it is likely that subtrees taken from existing, partially evolved models will be more useful than new, randomly generated subtrees, crossover is applied more frequently than mutation. This process is iterated over many evolutionary generations, each time replacing the population with a new population of altered trees. Over time, this produces populations of increasing fitness.

The average wall-clock time for one experiment using the Vermont Advanced Computing Core (VACC) supercomputer was 333 s for Experiment Set I (3000 generations) and 1207 s for Experiment Set II (10,000 generations). The total wall-clock time for all of Experiment Set I was approximately 89 h. The total wall-clock time for all of Experiment Set II was approximately 321 h. Because GP is a stochastic optimization method, its computation complexity is unclear. However, recent work has begun to address this problem (Neumann et al., 2011; Durrett et al., 2010).

One challenge facing GP, like all techniques for deriving a model from training data, is over-fitting. An over-fit model performs well on training data but does not generalize well and fails on unseen data. It memorizes values instead of capturing the mathematical relationships among the data.

The size of a GP model (number of nodes in a tree) constrains its complexity and fitness. Trees that are too small are too simple to accurately model the data and are under-fit. They perform poorly on both training and testing data. Trees that become too large perform extremely well on training data but, due to over-fitting, perform poorly on unseen data. Somewhere between these extremes lies the best, non-over-fit model.

In order to explore the gradient from small, under-fit models to large, over-fit models, we added multi-objective Pareto optimization to lil-gp. Pareto optimization applies evolutionary pressure toward multiple simultaneous goals, in this case low error and small model size, by producing a population (front) of nondominated models. A tree is dominated by another tree if it is inferior by all objectives, i.e. it is both larger and has lower fitness. A Pareto front (non-dominated front) consists of a set of trees such that no tree is dominated by any other tree on the front. The non-dominated trees are selected at each GP generation so that each population is a non-dominated front, including the final population. The result of GP is therefore a set of trees of various sizes. We set an absolute upper bound at size 30 because we had observed that models with size larger than 30 were consistently over-fit. Arranged from smallest to largest, the error of these trees on the training data decreases monotonically. Error on unseen data, however, will decrease only to a point, and will then increase beyond some tree size as the models become over-fitted.

At this point is the tree size that will maximize performance on ϱ without over-fitting. Models no bigger than this can express

features common to both training and testing data but cannot express features that are unique to the training data. However, this size threshold is not known while generating models because test data are not available. It must remain *unseen* for model testing. We therefore developed a novel *selection set* method for selecting a single model from the Pareto front. In the *selection set* method, the training data are further divided into two subsets of equal size, a growth set, g, and a selection set, s (Eq. 2). GP is applied to g to obtain a Pareto front. Each model on the front is then evaluated on s. GP returns the model that performs best (lowest error) on s. We used the *election set* method in all experiments.

4.6. Binary regression trees

We include BT in Experiment Set II in order to compare GP to another nonlinear, less computationally demanding, modeling technique. Erxleben et al. (2002) compared the performances of four spatial interpolation methods to estimate SWE and found that a method combining binary regression trees with geostatistical methods was more accurate than other methods. We used the DecisionTreeRegressor class of the Scikit-learn machine learning module for Python (Pedregosa et al., 2011). This software implements the Classification and Regression Trees (CART) algorithm, which is similar to C4.5 (Hastie et al., 2009). BT is parameterized by the maximum tree depth; we used default options for other parameters. As with GP, the data for BT was divided into g, s, and τ . For each experiment, a set of trees was trained on g such that the *n*th tree had a maximum depth of *n*. The maximum value of *n* was determined by incrementing *n* until further increase did not result in larger trees. The maximum value of *n* varied between 7 and 13.

Like the Pareto front produced by GP with multi-objective optimization, this methods results in a gradient of models ranging from very small models with high error on g to very large models with low error on g. Each is evaluated on s and the one with the lowest error is returned by BT to be evaluated on τ in order to determine model error. Thus, we applied the same *selection set* method to BT as to GP in order to discourage over-fitting and to provide similar exposure to the data so that the performance of the techniques may be compared. Note, however, that in the case of GP, multi-objective optimization applies pressure toward model parsimony continuously over the course of the evolution of a population of models. In the case of BT, the selection set method was applied once to a set of models after they have been generated.

5. Experiments: descriptions and results

In this section we describe the experiments we conducted and report the results.

5.1. Experiment Set I

In Experiment Set I measurements from snow courses provided ground-truth *SWE* data. We developed models to predict snow course *SWE* at eight different sites in California where snow courses had been conducted (Table 1). Three sites (CAP, GRZ, KTL) are located at snow pillows but are not near any NCDC weather stations. Three sites (NTH, SPD, MSH) are near NCDC stations but are not at snow pillows. Two of the snow course sites (HYS and HIG) are located at snow pillows and are also near NCDC stations.

First, we conducted experiments at sites with snow pillows but without weather stations (CAP, GRZ, KTL). These experiments explored how well linear and nonlinear models predict snow

Table 5Experiment Set I summary.

Experiment	Model inputs	Locations
a	air temp.	MSH, NTH, SPD, HIG, HYS
b	TOY	all
c	pillow	CAP, GRZ, KTL, HIG, HYS
d	air temp., TOY	MSH, NTH, SPD, HIG, HYS
e	air temp., pillow	HIG, HYS
f	TOY, pillow	CAP, GRZ, KTL, HIG, HYS
g	air temp. TOY pillow	HTC HYS

course-derived ground truth *SWE* using only snow pillow measurements. Inputs to the models were pillow *SWE* and *TOY*. At each site we developed models with three combinations of input variables: *TOY* alone, pillow *SWE* alone, and *TOY* combined with pillow *SWE*. In each case, we compared the performance of GP, LR, and BM.

Second, we conducted experiments at sites near weather stations but without snow pillows (\mathcal{KTL} , \mathcal{MSH} , \mathcal{NTH}). These experiments explored how well linear and nonlinear models predict snow course-derived ground truth *SWE* using air temperature data without access to snow pillow *SWE* measurements. Inputs to the models were *air temperature* and *TOY*. At each site we develop models with three combinations of input variables: temperature alone, *TOY* alone, and temperature combined with *TOY*. In each case, we compare the performance of GP to LR. Third, we conducted experiments at sites that are near weather stations and have snow pillows (\mathcal{HIG} , HYS). These experiments explored how well linear and nonlinear models predict snow course-derived ground truth *SWE* using both pillow *SWE* measurements and air temperature data. Inputs to the models were *SWE*, *air temperature*, and *TOY*. At each site we develop models with seven unique combinations of input variables: temperature alone, *TOY* alone, pillow *SWE* alone, temperature and *TOY* together, temperature and pillow *SWE* together, *TOY* and pillow *SWE* together, and, finally, temperature, *TOY*, and pillow *SWE* together.

Table 5 summarizes Experiment Set I. Each experiment was repeated 30 times to generate error samples for each method. Figs. 6–9 plot the mean values of the samples. Error bars indicate 95% confidence intervals, i.e. sample mean \pm (SEM × 1.96). GP and LR had similar error, but both had lower error than BM with *p*-value less than 0.001 in all cases.

The mean ground truth *SWE* value in mm at each site was: CAP : 1145, GRZ : 1256, KTL : 687, MSH : 1747, NTH : 337, SPD : 697, HIG : 594, HYS : 1065..

5.2. Experiment Set II

In Experiment Set II models predicted *HS* instead of *SWE*. While research on the influence of meteorological factors on snowpack distribution is extensive (Logan, 1973; Elder et al., 1991;



Fig. 6. Experiment Set I results: CAP, GRZ, and KTL.



Fig. 7. Experiment Set I results: MSH, NTH, and SPD.



Fig. 9. Experiment Set I results: HYS.

Schmucki et al., 2014; Hock and Noetzli, 1997), the inclusion of meteorological inputs does not always improve snowpack model performance (Moeser, 2010), and the inclusion of air temperature data did not improve model performance in Experiment Set I. Therefore, in Experiment Set II we focus on *TOY* and single-point *HS* measurements as predictors of mean catchment *HS*. Instead of manual snow course data as in Experiment Set I, ground-truth data are derived from *HS* measurements collected by the Snowcloud WSN. We compared the performance of three machine learning techniques: LR, BT, and GP.

We developed estimators to predict *HS* at two sites: Sulitjelma, Norway and the Sagehen Experimental Forest, California. At Sulitjelma, model inputs were combinations of *HS* at Balvatn and *TOY*. At Sagehen, model inputs were combinations of *HS* at \mathcal{HYS} , *HS* at \mathcal{TDC} , and *TOY*. Table 6 summarizes Experiment Set II. We repeated each experiment four times (*Random Division*, 4 *Bins*, 3 *Bins*, 2 *Bins*) and each of these 30 times to generate error samples.

Figs. 10–13 plot the mean values of the samples, i.e. the error of the modeling techniques on testing data. Error bars indicate 95% confidence intervals, i.e. sample mean \pm (SEM \times 1.96). Stars indicate *p*-values for the Student's paired *t*-test with the hypothesis

Table 6	
Experiment Set II summary.	

Experiment	Location	Model inputs
a	Sulitjelma, Norway	TOY
b	Sulitjelma, Norway	HS at Balvatn
с	Sulitjelma, Norway	HS at Balvatn, TOY
d	Sagehen, California	TOY
e	Sagehen, California	HS at HYS
f	Sagehen, California	HS at <i>IDC</i>
g	Sagehen, California	HS at HYS, TOY
h	Sagehen, California	HS at <i>IDC</i> , TOY



Fig. 10. Experiment Set II (random division) model error.



Fig. 11. Experiment Set II (four bins) model error.

the GP does not have lower error than BT, i.e. the probability that GP does not outperform BT. One star, *, indicates that p is less than 0.05, ** indicates that p is less than 0.01, and *** indicates that p is less than 0.001. Similarly, plus signs indicate p-values for the hypothesis that GP does not have lower error than LR, i.e. the probability that GP does not outperform LR. One plus sign, +, indicates that p is less than 0.05, and ++ indicates that p is less than 0.01. The mean ground truth HS value at Sulitjelma was 1.1900 m. The mean ground truth HS value at Sagehen was 0.728 m.

Figs. 14–17 plot the mean sizes of the models whose performance is reported in Figs. 10–13. In the case of GP and BT, these



Fig. 12. Experiment Set II (three bins) model error.



Fig. 13. Experiment Set II (two bins) model error.



Fig. 14. Experiment Set II (random division) model size.

are the models selected using the *selection set* method. For GP, model size is the number of nodes in the GP tree. For BT, model size is the number of nodes in the binary tree. For LR, model size is the number of operators and values, specifically 5 in the case of a single independent variable and 9 in the case of two independent variables. Stars indicate *p*-values for the Student's paired *t*-test with the hypothesis the GP models are not smaller than BT models. One star, *, indicates that *p* is less than 0.05, ** indicates that *p* is less than 0.001.

6. Discussion

In this section we discuss the results of our experiments, offer some hypotheses to explain our findings, and suggest possible next steps for continued research.

6.1. Experiment Set I

In Experiment Set I GP performed at least as well as other methods in all experiments. This result was expected because GP is capable of generating the same models as LR and BM. We did not perform hypothesis tests comparing GP with LR because visual inspection of error means and 95% confidence intervals (Figs. 6– 9) suggests that the methods performed similarly. At the sites where a snow pillow was present, the performance of BM was evaluated. At all of these sites, in all of the experiments where pillow *SWE* was an input variable (b, c, f), both LR and GP performed better (*p*-value less than 0.001) than BM.

These results suggest that machine learning techniques can be used to develop models that predict mean catchment *SWE* more accurately than BM. In general, models performed better when





Fig. 16. Experiment Set II (three bins) model size.



Fig. 17. Experiment Set II (two bins) model size.

snow pillow data were included. However, GP did not outperform LR.

Because LR performed as well as GP in Experiment Set I, we suspected strict linearity among the explanatory relationships in the data. We hypothesize that because snow courses measure *SWE* only at a single location, they failed to capture existing nonlinearities, and that even though the relationships underlying snowpack distribution are nonlinear, our Experiment Set 1 data is linear. We therefore did not further pursue nonlinear modeling, such as BT, in Experiment Set 1.

6.2. Experiment Set II

First we conducted Experiment Set II: *Random Division*. GP outperformed LR in every experiment except in Norway when the only model input was HS at Balvatn. In every experiment in California where TOY was an input, BT has much lower error than either GP or LR. In all experiments where TOY was an input, the resulting BT models were very large. GP also had lower error and larger model sizes when TOY was used then when TOY was not used. We had originally introduced the TOY variable to allow models to distinguish different parts of the season. However, we hypothesized the BT, and to a lesser extent GP, were abusing the TOY variable to memorize snow data by mapping TOY data to ground truth HS. Even though training and testing data were technically distinct, many of the samples in the testing data were temporally or spatially proximal to samples in the training data. The testing data were not truly unseen with respect to the TOY variable. Even though models generalized well to the testing data, they were over-fitting to the TOY variable and would likely not generalize to truly unseen data, e.g. from another snow season.

To test this hypothesis and address the possible problem of over-fitting to the *TOY* variable, we repeated Experiment Set II three more times. In Experiment Set II: *4 Bins*, *3 Bins*, and *2 Bins*, we successively decreased the temporal overlap between training and testing data and increase the coarseness of the temporal granularity of the division into training and testing data. Proceeding through this sequence, it became more difficult for machine learning to memorize *HS* data by over-fitting to the *TOY* variable. At the same time, BT error increased and the performance of GP with respect to BT improved. These results suggest that GP is more resilient against over-fitting than BT, possible as a result of multi-objective optimization. Furthermore, when the ability of machine learning to exploit the *TOY* variable by memorizing *HS* the data were minimized, GP significantly outperformed both LR and BT.

6.3. Future work

We believe that the preliminary results discussed in this work are promising and warrant further research into of the applicability of GP to snowpack modeling.

This work should be expanded into a multi-year study. Although Experiment I used snow course data collected over several years, Snowcloud data used in Experiment II was limited to single snow season. A multi-year study would allow models trained on Snowcloud data during one or several years to be evaluated on unseen data from another year. Models trained on multi-year data may be more robust to application in future years than are models trained on single-year data. Even without collecting more data, Experiment Set I could be modified so that models are trained on data from earlier years and tested on data from later years.

Beyond those discussed here, there are many machine learning techniques that should be applied to the problem of catchment-scale *SWE* estimation. GP possesses a unique combination of desirable qualities, but its performance should be compared against other methods such as ANNs, nonlinear multiple regression, and FFX (McConaghy, 2011), a non-evolutionary symbolic regression technology.

The only meteorological input to our models was air temperature. However, meteorological data involving wind, solar radiation, humidity, etc. are available for many locations and have been shown to influence snow distribution (Logan, 1973; Elder et al., 1991; Schmucki et al., 2014; Hock and Noetzli, 1997). Future work should incorporate more potential meteorological predictors of *SWE* and *HS*.

Topographic features significantly shape snow distribution, and models of this relationship have been developed and used extensively (Winstral et al., 2013; Marofi et al., 2011; Chang and Li, 2000; Tabari et al., 2010; Anderton et al., 2004; Grünewald et al., 2013; Molotch et al., 2005; Elder et al., 1998). Although topographic data was not an explicit input in our experiments, models developed with our techniques that use input variables to predict distributed snow measurements likely express some of the relationship between topography and snowpack distribution. Previous efforts to model snowpack using topographic data have derived explicit model inputs from DEMs. The possibility that GP could play an active role in determining which topographical features to use should be explored. GP might discover new methods for extracting information from DEMs that is predictive of snowpack distribution. It is possible that machine learning could use topographic and other data to produce non-cite-specific models, which are trained on data from one or more site and then applied to other sites.

Schwaerzel and Bylander (2006) developed high-order statistical functions for GP to model financial data. These allowed GP models to dynamically select and aggregate a slice of time series data. Future work should apply these techniques to allow GP to determine how to select and aggregate meteorological and topographic data. We made air temperature available to GP by means of functions that aggregate daily measurements over an arbitrary seven day window. Instead, GP could inductively discover how models should dynamically select and aggregate a section of time series data according to changing circumstances.

7. Conclusion

In this paper we have described novel, low-cost methods for catchment-scale SWE estimation using machine learning algorithms. The commonly used method of estimating catchment-scale SWE from a single point measurement is error-prone because of the spatial heterogeneity of snowpack distribution. We envision an approach wherein short-term field campaigns collect groundtruth data for generating snowpack models which can subsequently augment existing NRT snow telemetry. Toward this end, we explored a suite of machine learning techniques to extrapolate estimates of mean catchment SWE from single point SWE measurements and other available data and pursued three key research directions. First, we addressed the question of which machine learning approaches are best for this problem. Second, we discussed and pursued the use of a range of possible input parameters. Finally, we grappled with the issue of ground-truthing given limited datasets.

We compared the performance of a basic method (BM) which assumes no spatial variability of SWE, linear regression (LR), Genetic Programming (GP), and binary regression trees (BT). We emphasize GP because it produces nonlinear, white-box models without requiring assumptions about model form. GP can be augmented with multi-objective optimization to constrain model complexity and mitigate over-fitting. We found that machine learning techniques generally outperformed BM, demonstrating the spatial variability of SWE. Nonlinear techniques outperformed linear models in Experiment Set II, but not in Experiment Set I, suggesting that there are nonlinear relationships among the modeled data used in Experiment Set II. Snowpack distribution at the catchment scale has been shown to be highly nonlinear. It is possible that the spatially distributed sampling technique (Snowcloud wireless sensor network) used for ground-truthing in Experiment Set II captured some of the nonlinearity of snowpack distribution, while the single-location sampling (manual snow courses) used for Experiment Set I did not.

When we naively divided our data at random to generate training and testing data, BT had much lower error than GP in experiments where time of year (*TOY*) was an input variable. In these cases, BT models were much larger than PG models and we suspected that they were memorizing data by mapping *TOY* to snow depth. When we instead divided the data into more temporally contiguous training and testing data in order to prevent this behavior, BT model size decreased and GP outperformed BT.

We emphasize that GP can flexibly incorporate new predictors of catchment-scale *SWE* into the models generated, augmenting its capacity to extrapolate estimates of mean catchment-wide *SWE* from a single point measurement. Genetic programming will make use of input data that helps explain the dependent variable while ignoring data that does not. Our choice of independent variables was a result of intuitive guesses combined with constraints on available data. Topographic information was ruled out because we were unable to determine the precise locations of snow pillows. Multiple forms of meteorological data were available, but air temperature was the most complete, allowing us to compose datasets large enough for effective machine learning. However, the inclusion of air temperature did not have a significant impact on model performance in our first experiment set, and so we did not use any meteorological data in our second experiment set.

Because it has been shown that the forcing effects underlying snowpack distribution change over the course of a snow season, we introduced time of year (*TOY*) as an independent variable so that models can distinguish seasonal differences. However, we found that nonlinear models used *TOY* to memorize the data by mapping *TOY* to ground truth measurements instead of expressing the underlying relationships of snowpack distribution. The ideal solution to this problem would be a multi-year study using spatially distributed data collected by Snowcloud. However, given the limitation of a one year dataset, we modified how data was divided to constrain the temporal proximity of training and testing data.

We conducted two sets of experiments, using available data, as successive approximations of our goal of near-real-time catchment-scale *SWE* estimation. When ground truth was obtained from distributed sampling techniques and when we were careful to mitigate overfitting to the *TOY* variable, GP outperformed other techniques.

Acknowledgments

We would like to acknowledge several individuals for contributions to the content of this paper. Dr. Ian Brown, Stockholm University. Dr. Jeff Frolik, University of Vermont. Dr. Jeff Dozier, University of California. Jeff Brown, University of California. David Moeser, WSL-Institut für Schnee- und Lawinenforschung SLF in Davos, Switzerland. Dr. Keith Klepeis, University of Vermont. Rune Engeset, Norwegian Water Resources and Energy Directorate. Heidi Bache Stranden, Norwegian Water and Energy Directorate.

We acknowledge the Vermont Advanced Computing Core which is supported by NASA (NNX 06AC88G), at the University of Vermont for providing High Performance Computing resources that have contributed to the research results reported within this paper. We acknowledge the support from DARPA through grants W911NF-11-1-0076 and FA8650-11-1-7155. We acknowledge the support from NSF through grant #PECASE-0953837. We acknowledge the support of the Air Force Office of Scientific Research through a YIP grant. This work was supported by NASA under Cooperative Agreement #NNX10AK67H-S02.

References

- Adams, W.P., 1976. Areal differentiation of snow cover in east central Ontario. Water Resour. Res. 12, 1226–1234.
- Anderton, S., White, S., Alvera, B., 2004. Evaluation of spatial variability in snow water equivalent for a high mountain catchment. Hydrol. Process. 18, 435–453.
 Bales, R.C., Molotch, N.P., Painter, T.H., Dettinger, M.D., Rice, R., Dozier, J., 2006.
- Mountain hydrology of the western united states. Water Resour. Res., 42 Balk, B., Elder, K., 2000. Combining binary decision tree and geostatistical methods
- to estimate snow distribution in a mountain watershed. Water Resour. Res. 36, 13–26.

- Boxalla, B., 2014. California Snowpack Hits Record Low. http://articles.latimes.com/2014/jan/30/local/la-me-brown-water-20140131>.
- Bühler, Y., Christen, M., Kowalski, J., Bartelt, P., 2011. Sensitivity of snow avalanche simulations to digital elevation model quality and resolution. Ann. Glaciol. 52, 72–80.
- Chang, K.T., Li, Z., 2000. Modelling snow accumulation with geographic information system. Int. J. Geogr. Inf. Sci. 14, 693–707.
- Dozier, J., 2011. Mountain hydrology, snow color, and the fourth paradigm. Eos Trans. Am. Geophys. Union 92, 373–374.
- Dozier, J., Painter, T.H., 2004. Multispectral and hyperspectral remote sensing of alpine snow properties. Annu. Rev. Earth Planet. Sci. 32, 465–494.
- Durrett, G., Neumann, F., O'Reilly, U., 2010. Computational complexity analysis of simple genetic programming on two problems modeling isolated program semantics. CoRR abs/1007.4636.
- Elder, K., Dozier, J., Michaelsen, J., 1991. Snow accumulation and distribution in an alpine watershed. Water Resour. Res. 27, 1541–1552.
- Elder, K., Rosenthal, W., Davis, R.E., 1998. Estimating the spatial distribution of snow water equivalence in a montane watershed. Hydrol. Process. 12, 1793–1808.
- Engeset, R., Tveito, O.E., Alfnes, E., Mengistu, Z., Udns, H.C., Isaksen, K., Frland, E.J., 2004. Snow map system for Norway, in: Proc. Nordic Hydrol. Conf., p. 12.
- Erxleben, J., Elder, K., Davis, R., 2002. Comparison of spatial interpolation methods for estimating snow distribution in the colorado rocky mountains. Hydrol. Process. 16, 3627–3649.
- Fassnacht, S., Dressler, K., Bales, R., 2003. Snow water equivalent interpolation for the colorado river basin from snow telemetry (snotel) data. Water Resour. Res., 39
- lil-gp Genetic Programming System, 2013. http://garage.cse.msu.edu/software/lil-gp/
- Grünewald, T., Stotter, J., Pomeroy, J., Dadic, R., Baños, I.M., Marturià, J., Spross, M., Hopkinson, C., Burlando, P., Lehning, M., 2013. Statistical modelling of the snow depth distribution in open alpine terrain. Hydrol. Earth Syst. Sc., 17
- GSL, 2014. GNU Scientific Library. < http://www.gnu.org/software/gsl/>
- Guan, B., Molotch, N.P., Waliser, D.E., Fetzer, E.J., Neiman, P.J., 2010. Extreme snowfall events linked to atmospheric rivers and surface air temperature via satellite measurements. Geophys. Res. Lett., 37
- Hastie, T., Tibshirani, R., Friedman, J., Hastie, T., Friedman, J., Tibshirani, R., 2009. The Elements of Statistical Learning, vol. 2. Springer.
- Hock, R., Noetzli, C., 1997. Areal melt and discharge modelling of Storglaciären. Sweden. Ann. Glaciol. 24, 211–217.
- Johnson, J.B., Marks, D., 2004. The detection and correction of snow water equivalent pressure sensor errors. Hydrol. Process. 18, 3513–3525.
- Jost, G., Weiler, M., Gluns, D.R., Alila, Y., 2007. The influence of forest and topography on snow accumulation and melt at the watershed-scale. J. Hydrol. 347, 101–115.
- Kerkez, B., Glaser, S.D., Bales, R.C., Meadows, M.W., 2012. Design and performance of a wireless sensor network for catchment-scale snow and soil moisture measurements. Water Resour. Res., 48
- Koza, J.R., 1992. Genetic Programming. Massachusetts Institue of Technology, Cambridge, MA.
- Logan, L., 1973. Basin-wide water equivalent estimation from snowpack depth measurements. Role Snow Ice Hydrol, vol. 107. IAHS AIHS Publ., pp. 864–884.
- López-Moreno, J., Fassnacht, S., Heath, J., Musselman, K., Revuelto, J., Latron, J., Morán-Tejeda, E., Jonas, T., 2012. Small scale spatial variability of snow density and depth over complex alpine terrain: implications for estimating snow water equivalent. Adv. Water Resour.
- Marks, D., Domingo, J., Susong, D., Link, T., Garen, D., 1999. A spatially distributed energy balance snowmelt model for application in mountain basins. Hydrol. Process. 13, 1935–1959.
- Marofi, S., Tabari, H., Abyaneh, H.Z., 2011. Predicting spatial distribution of snow water equivalent using multivariate non-linear regression and computational intelligence methods. Water Resour. Manag. 25, 1417–1435.
- Martinec, J., Rango, A., 1981. Areal distribution of snow water equivalent evaluated by snow cover monitoring. Water Resour. Res. 17, 1480–1488.
- McConaghy, T., 2011. FFX: Fast, scalable, deterministic symbolic regression technology. In: Genetic Programming Theory and Practice IX. Springer, pp. 235–260.
- Meromy, L., Molotch, N.P., Link, T.E., Fassnacht, S.R., Rice, R., 2013. Subgrid variability of snow water equivalent at operational snow stations in the western USA. Hydrol. Process. 27, 2383–2400.
- Milly, P., Betancourt, J., Falkenmark, M., Hirsch, R., Kundzewicz, Z., Lettenmaier, D., Stouffer, R., 2008. stationarity is dead: whither water management? Science 319, 573–574.
- Moeser, C.D., 2010. Development, Analysis and Use of a Distributed Wireless Sensor Network for Quantifying Spatial Trends of Snow Depth and Snow Water Equivalence Around Meteorological Stations With and Without Snow Sensing Equipment. Master's thesis. University of Nevada – Reno.
- Moeser, C.D., Walker, M., Skalka, C., Frolik, J., 2011. Application of a wireless sensor network for distributed snow water equivalence estimation, in: Proc. West. Snow Conf., Stateline, NV, USA.
- Molotch, N., Colee, M., Bales, R., Dozier, J., 2005. Estimating the spatial distribution of snow water equivalent in an alpine basin using binary regression tree models: the impact of digital elevation data and independent variable selection. Hydrol. Process. 19, 1459–1479.
- Molotch, N.P., Bales, R.C., 2005. Scaling snow observations from the point to the grid element: implications for observation network design. Water Resour. Res., 41.

- Molotch, N.P., Bales, R.C., 2006. Snotel representativeness in the rio grande headwaters on the basis of physiographics and remotely sensed snow cover persistence. Hydrol. Process. 20, 723–739.
- NASA Airborne Snow Observatory, 2015. Measuring Spatial Distribution of Snow Water Equivalent and Snow Albedo. http://aso.jpl.nasa.gov/.
- National Snow & Ice Data Center, 2014. Clpx overview. http://nsidc.org/data/clpx/. Neumann, F., O'Reilly, U.M., Wagner, M., 2011. Computational complexity analysis of genetic programming – initial results and future directions. In: Riolo, R., Vladislavleva, E., Moore, J.H. (Eds.), Genetic Programming Theory and Practice IX. Springer New York, Genetic and Evolutionary Computation, pp. 113–128.
- Ohmura, A., 2001. Physical basis for the temperature-based melt-index method. J. Appl. Meteorol. 40, 753–761.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E., 2011. Scikit-learn: machine learning in Python. J. Mach. Learn. Res. 12, 2825–2830.
- Pierce, D.W., Barnett, T.P., Hidalgo, H.G., Das, T., Bonfils, C., Santer, B.D., Bala, G., Dettinger, M.D., Cayan, D.R., Mirin, A., et al., 2008. Attribution of declining western US snowpack to human effects. J. Clim., 21
- Rittger, K., Kahl, A., Dozier, J., 2011. Topographic distribution of snow water equivalent in the sierra nevada, in: Proc. West. Snow Conf., Western Snow Conference.
- Rittger, K.E., 2012. Spatial Estimates of Snow Water Equivalent in the Sierra Nevada. Ph.D. Thesis. University Of California Santa Barbara.
- Schirmer, M., Wirz, V., Clifton, A., Lehning, M., 2011. Persistence in intra-annual snow depth distribution: 1. Measurements and topographic control. Water Resour. Res. 47, W09516.
- Schmidt, M.D., Vallabhajosyula, R.R., Jenkins, J.W., Hood, J.E., Soni, A.S., Wikswo, J.P., Lipson, H., 2011. Automated refinement and inference of analytical models for metabolic networks. Phys. Biol. 8, 055011.
- Schmucki, E., Marty, C., Fierz, C., Lehning, M., 2014. Evaluation of modelled snow depth and snow water equivalent at three contrasting sites in Switzerland using SNOWPACK simulations driven by different meteorological data input. Cold Reg. Sci. Technol. 99, 27–37.
- Schwaerzel, R., Bylander, T., 2006. Predicting Financial Time Series by Genetic Programming with Trigonometric Functions and High-Order Statistics, GECCO.

- Scipión, D., Mott, R., Lehning, M., Schneebeli, M., Berne, A., 2013. Seasonal smallscale spatial variability in alpine snowfall and snow accumulation. Water Resour. Res. 49, 1446–1457.
- Serreze, M.C., Clark, M.P., Armstrong, R.L., McGinnis, D.A., Pulwarty, R.S., 1999. Characteristics of the western united states snowpack from snowpack telemetry (snotel) data. Water Resour. Res. 35, 2145–2160.
- Skalka, C., Frolik, J., 2014. Snowcloud: a complete data gathering system for snow hydrology research. In: Real-World Wireless Sensor Networks. Springer, pp. 3–14.
- Snow Surveyor, 2014. http://www.water.ca.gov/floodmgmt/hafoo/hb/sss/surveyor.cfm.
- Sturm, M., Taras, B., Liston, G.E., Derksen, C., Jonas, T., Lea, J., 2010. Estimating snow water equivalent using snow depth data and climate classes. J. Hydrometeorol., 11.
- Tabari, H., Marofi, S., Abyaneh, H.Z., Sharifi, M.R., 2010. Comparison of artificial neural network and combined models in estimating spatial distribution of snow depth and snow water equivalent in Samsami basin of Iran. Neural Comput. Appl. 19, 625–635.
- Tappeiner, U., Tappeiner, G., Aschenwald, J., Tasser, E., Ostendorf, B., 2001. GISbased modelling of spatial pattern of snow cover duration in an alpine area. Ecol. Model. 138, 265–275.
- USDA, 2014. Snow Surveys and Water Supply Forecasting. http://www.nrcs.usda.gov/wps/portal/nrcs/detail/or/snow/?cid=nrcs142p2_046152.
- Welch, S.C., Kerkez, B., Bales, R.C., Glaser, S.D., Rittger, K., Rice, R.R., 2013. Sensor placement strategies for snow water equivalent (SWE) estimation in the american river basin. Water Resour. Res. 49, 891–903.
- Winstral, A., Elder, K., Davis, R.E., 2002. Spatial snow modeling of windredistributed snow using terrain-based parameters. J. Hydrometeorol. 3, 524– 538.
- Winstral, A., Marks, D., 2014. Long-term snow distribution observations in a mountain catchment: assessing variability, time stability, and the representativeness of an index site. Water Resour. Res. 50, 293–305.
- Winstral, A., Marks, D., Gurney, R., 2009. An efficient method for distributing wind speeds over heterogeneous terrain. Hydrol. Process. 23, 2526–2535.
- Winstral, A., Marks, D., Gurney, R., 2013. Simulating wind-affected snow accumulations at catchment to basin scales. Water Resour. Res. 55, 64–79.