

OPTIMIZATION METHODS FOR MACHINE LEARNING

BETHANY LUSCH Asst. Computer Scientist Argonne National Lab Leadership Computing Facility blusch@anl.gov



August 9, 2019 ATPESC



- "best route" (Minimize cost of delivery subject to all mail is delivered)
- "best product"
- (Minimize -profit subject to safety)
- "best prediction model" (Minimize prediction error)





MACHINE LEARNING

$\begin{array}{ll} \underset{x}{\text{minimize}} & f(x) \\ \text{subject to} & \dots \text{constraints...} \end{array}$

- "best prediction model"
- "best recommendation"
- "best clusters"

(Minimize prediction error)(Minimize # who don't buy anything)(Minimize distance within clusters while maximizing distance between clusters)

Underneath most ML problems is an optimization problem



TYPES OF OPTIMIZATION

- Linear
- Quadratic
- Convex

General

 $\begin{array}{ll} \underset{x}{\text{minimize}} & f_0(x)\\ \text{subject to} & f_i(x) \leq b_i, \ i = 1, \dots, m. \end{array}$

- Have 2nd derivs
- Have gradients



Desperation



TYPES OF OPTIMIZATION

- Linear
- Quadratic
- Convex
- Have 2nd derivs
- Have gradients

General

U.S. DEPARTMENT OF ENERGY Argonne National Laboratory is a U.S. Department of Energy laboratory managed by UChicago Argonne, LLC Roughly...

 More time on formulating problem to fit these categories

More time on optimization algorithm











General

U.S. DEPARTMENT OF ENERGY Argonne National Laboratory is a U.S. Department of Energy laboratory managed by UChicago Argonne, LLC If Q is positive definite: Weakly polynomial-time algorithms Local minima are global optima

CONVEX OPTIMIZATION



DIFFERENTIABLE OPTIMIZATION



9

GENERAL OPTIMIZATION

- Linear
- Quadratic
- Convex

Desperation

- Have 2nd derivs
- Have gradients





minimize

Generally NP-hard Hopefully know *some* structure!



 $f_0(x)$

subject to $f_i(x) \leq b_i, i = 1, \dots, m$.

DISCRETE OPTIMIZATION

- Linear
- Quadratic
- Convex

Desperation

- Have 2nd derivs
- Have gradients







Picture source: wallpaperflare.com

CLASSIFICATION EXAMPLE

- Problem: label each document x as related to politics or not (1 or -1).
- Hard to come up with rules by hand, so ML helps: learn function h(x)
- Really want to minimize expected risk of misclassification:

$$\underset{h}{\text{minimize}} \quad R(h) = \mathbb{P}[h(x) \neq y] = \mathbb{E}[\mathbb{1}[h(x) \neq y]]$$

- How do we pick family of functions to optimize over?
- How do we know which one is optimal?





REALITIES

Really want to minimize expected risk of misclassification:

$$\underset{h}{\text{minimize}} \quad R(h) = \mathbb{P}[h(x) \neq y] = \mathbb{E}[\mathbb{1}[h(x) \neq y]]$$

Don't know probability distribution, so minimize empirical risk:

minimize
$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(x_i) \neq y_i]$$

Easier if smooth loss and parameterized h:

$$\underset{w}{\text{minimize}} \quad R_n(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w), y_i)$$



CHOOSING FUNCTION FAMILY

- Possibility of low empirical risk on training data
- Expected risk and empirical risk don't have large gap
- Can efficiently solve optimization
- (convenient representation, smoothness,...)

Bias vs. variance Overfitting vs. underfitting

Major themes of machine learning!



BIAS VS. VARIANCE





Picture similar to: Seema Singh, "Understanding the Bias-Variance Tradeoff"



LINEAR REGRESSION (LEAST-SQUARES)

Desperation



- Quadratic
- Convex
- Have 2nd derivs
- Have gradients

General

 $\underset{x}{\operatorname{minimize}} \|Ax - b\|_2^2$

if multiply out: $\begin{array}{ll} \underset{x}{\text{minimize}} & x^{t}A^{T}Ax - 2b^{T}Ax + b^{T}b \\ \\ \underset{x}{\text{minimize}} & \underset{x}{\text{quadratic program!}} \end{array}$

Argonne



SUPPORT VECTOR MACHINE

Find linear classifier with maximum margin

- LinearQuadratic
- Convex

Desperation

- Have 2nd derivs
- Have gradients

General

U.S. DEPARTMENT OF ENERGY U.S. Department of Energy laborator managed by UChicago Argonne, LLC



SUPPORT VECTOR MACHINE

Find linear classifier with maximum margin

- LinearQuadratic
- Convex
- Have 2nd derivs
- Have gradients

General

Kernel SVM can do nonlinear classification while remaining a quadratic program

Argonne National Laboratory is a U.S. Department of Energy laborator managed by UChicago Argonne, LL

Desperation





-2

K-MEANS CLUSTERING

Desperation

Find clustering that minimizes distances within clusters



Picture source: Prasad Patil, "K Means Clustering : Identifying F.R.I.E.N.D.S in the World of Strangers'



RECALL: TYPES OF OPTIMIZATION

0
σ
<u> </u>
Φ
Q
S
Φ

 \cap

Linear

- Quadratic
- Convex
- Have 2nd derivs
- Have gradients

General

Roughly...

More time on formulating problem to fit these categories

More time on optimization algorithm





ANALOGOUSLY...

- Linear
- Quadratic
- Convex
- Have 2nd derivs
- Have gradients

General

Roughly...

More time on formulating problem (choosing features) so that these (biased) methods are suitable

More time on optimization algorithm



Desperation

REMINDER

- Possibility of low empirical risk on training data
- Expected risk and empirical risk don't have large gap
- Can efficiently solve optimization
- Big neural networks can be very expressive (low bias)
- So don't need to be as clever about input features
- But then easy to overfit...
- Optimization is tricky: optimization stalls, plus local minima or saddle points





Bias vs. variance Overfitting vs. underfitting





Picture source: Divakar Kapil in "Stochastic vs Batch Gradient Descent"





TYPES OF GRADIENT DESCENT

i.e. for empirical risk, explicitly summing over data points

$$R_n(w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$

 n_{i}

$$w_{k+1} \leftarrow w_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_{ik}(w_k)$$

n

Batch GD: use all points every step

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla f_{ik}(w_k)$$

Stochastic GD: use one point per step

$$w_{k+1} \leftarrow w_k - \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_{ik}(w_k)$$

Mini-batch GD: use a subset each step





TYPES OF GRADIENT DESCENT

Batch GD: use all points every step

Stochastic GD: use one point per step

Each step is accurate but expensive

Each step is noisy but fast

Mini-batch GD: use a subset each step

Happy medium?

Very common in deep learning, but often call it SGD





GRADIENT DESCENT CONSIDERATIONS

 $\underset{w}{\text{minimize}} \quad f(w)$

$$w_{k+1} \leftarrow w_k - \alpha_k \nabla f(w_k)$$

- Step size α_k
 - Too big: overshoot
 - Too small: very slow
 - (But can be good to escape local minima)
- Initialization
- Can you make the problem easier?



Li, et al. "Visualizing the Loss Landscape of Neural Nets" NeurIPS 2018





VARIANT: ADAM

Popular improvement on GD: Adam optimizer

- Separate learning rate for each weight
- Momentum: uses moving average of the gradient
- Also incorporates squared gradients

Cool exploration/visualization of momentum: https://distill.pub/2017/momentum/

(For those familiar: combines the best properties of AdaGrad, momentum, and RMSProp)





REGULARIZATION

- Common way to avoid overfitting: regularization
- Most common: L2 regularization





OVERFITTING CAUTION

- Can you generalize outside of your training set to validation/testing set?
- What about interpolating to data you haven't collected?
- Extrapolation extra unlikely to work





SUMMARY

- Linear
- Quadratic
- Convex

Desperation

- Have 2nd derivs
- Have gradients



Major themes in machine learning:

- Overfitting vs. underfitting
- Ability to efficiently solve optimization problem

For more, see:

SIAM REVIEW Vol. 60, No. 2, pp. 223–311	© 2018 Society for Industrial and Applied Mathematics

Optimization Methods for Large-Scale Machine Learning*

Léon Bottou[†] Frank E. Curtis[‡] Jorge Nocedal[§]



ANY QUESTIONS?

Thinking ahead to next talk: how would you parallelize gradient descent?





