

Machine Learning

Lecture 5 - Learning from large-scale data

Mattias Villani

Department of Statistics
Stockholm University

Department of Computer and Information Science
Linköping University



mattiasvillani.com



@matvil



mattiasvillani

Lecture overview

- More on loss functions
- Optimization algorithms for large data

Loss minimization as a proxy for generalization

- Parametric model $y|x \sim f_{\theta}(y|x)$.
- Learn the parameters by minimizing a cost function

$$\hat{\theta} = \arg \min_{\theta} \underbrace{J(\theta)}_{\text{cost function } J(\theta)} \equiv \underbrace{\frac{1}{n} \sum_{i=1}^n \overbrace{L(\hat{y}(x_i; \theta), y_i)}^{\text{loss function}}}_{\text{cost function } J(\theta)}$$

- ML: **cost function** $J(\theta)$ can be considered as a **proxy for** the real objective of interest, the **expected new data error**:

$$E_{\text{new}} \equiv \mathbb{E}_{\star} [E(\hat{y}(x_{\star}; \mathcal{T}), y_{\star})]$$

- **Early stopping** of iterative optimization:
 - ▶ when numerical accuracy is on par with statistical accuracy
 - ▶ to avoid overfitting (implicit regularization)

Loss functions for regression

■ Squared error loss

$$L(y, \hat{y}) = (\hat{y} - y)^2$$

■ Absolute error loss

$$L(y, \hat{y}) = |\hat{y} - y|$$

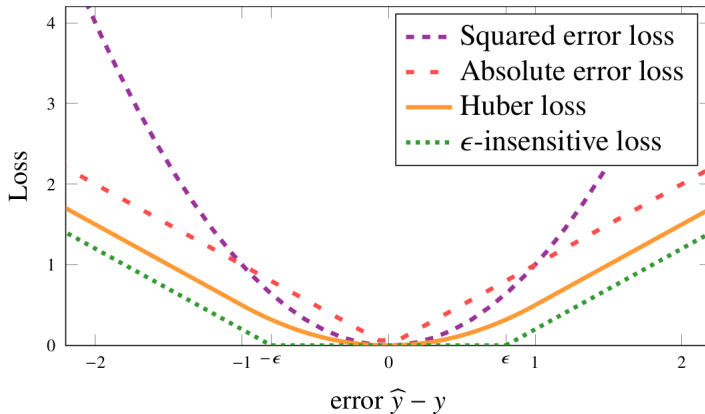
■ Huber loss

$$L(y, \hat{y}) = \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & \text{if } |\hat{y} - y| < 1 \\ |\hat{y} - y| - \frac{1}{2} & \text{otherwise} \end{cases}$$

■ ϵ -insensitive loss

$$L(y, \hat{y}) = \begin{cases} 0 & \text{if } |\hat{y} - y| < \epsilon \\ |\hat{y} - y| - \epsilon & \text{otherwise} \end{cases}$$

Loss functions for regression



Loss functions for classification

■ Misclassification loss

$$L(y, \hat{y}) = \mathbb{I}\{\hat{y} \neq y\} = \begin{cases} 0 & \text{if } \hat{y} = y \\ 1 & \text{if } \hat{y} \neq y \end{cases}$$

- ▶ doesn't care if a decision boundary is near a point or not.
- ▶ Gradient zero everywhere.

■ Better to use the probability in the loss: $\Pr(y = 1|x) = g(x)$.

■ Cross-entropy loss

$$L(y, g(x)) = \begin{cases} \log g(x) & \text{if } y = 1 \\ \log(1 - g(x)) & \text{if } y = -1 \end{cases}$$

■ Cross-entropy loss = Bernoulli distribution.

The margin for a classifier

- Class prediction by thresholding a function at zero:

$$\hat{y}(x) = \text{sign}(f(x))$$

- Example: logistic regression with $f(x) = x^\top \beta$.
- The **margin for a classifier** at (x, y) is

$$\text{margin} \equiv y \cdot f(x)$$

- Correct classification if and only if $f(x)$ and y have same sign.
- Misclassification loss in terms of margin

$$L(y \cdot f(x)) = \begin{cases} 1 & \text{if } y \cdot f(x) < 0 \\ 0 & \text{if } y \cdot f(x) > 0 \end{cases}$$

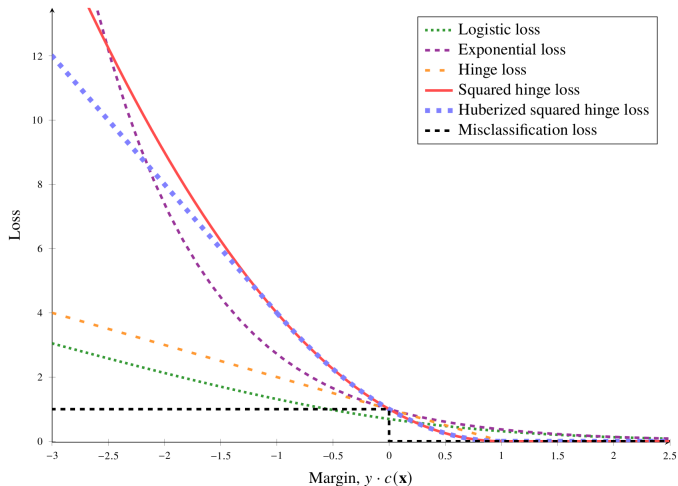
- The larger the margin, the more certain is the classifier.
- Log-likelihood for logistic regression (see eq. 3.34 in MLES):

$$L(y \cdot f(x)) = \log(1 + \exp(-y \cdot f(x)))$$

- **Exponential loss**

$$L(y \cdot f(x)) = \exp(-y \cdot f(x)).$$

Loss functions for classification



Log-likelihood as cost function

■ Negative log-likelihood as cost function

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n \log p(y_i | x_i; \theta)$$

■ Classification: cross-entropy.

■ Gaussian \implies squared loss. Laplace \implies absolute loss.

■ Advantages of the log-likelihood as loss function:

- ▶ ML estimator has good (asymptotic) properties.
- ▶ loss function comes from thinking about properties of the data: skewness, heavy-tails, truncation etc. We can **assess model fit**.
- ▶ natural extension to more complex problems:
 - censoring
 - missing data
 - dependent observations: $J(\theta) = -\log p(y_{1:n} | x_{1:n}; \theta)$.
- ▶ log prior for regularization in a Bayesian approach.
- ▶ log-likelihood is a **strictly proper** loss function.

Regularization revisited

- Note: MLES the loglikelihood is multiplied by the factor $\frac{1}{n}$

$$J(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \log p(y_i | \mathbf{x}_i; \boldsymbol{\theta})$$

- This changes the interpretation of λ in L2-regularization and the Ridge estimator

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

- **General explicit regularization**

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \lambda \cdot R(\boldsymbol{\theta})$$

- Implicit regularization:
 - ▶ Early stopping in iterative optimization methods
 - ▶ Dropout in neural networks.

Parameter optimization

■ Optimization problems in ML:

► Parameter learning

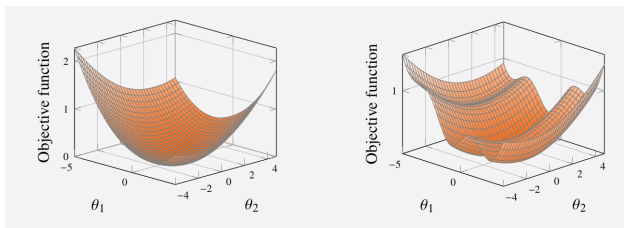
$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta; X, y)$$

$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta; X, y) + \lambda \cdot R(\theta)$$

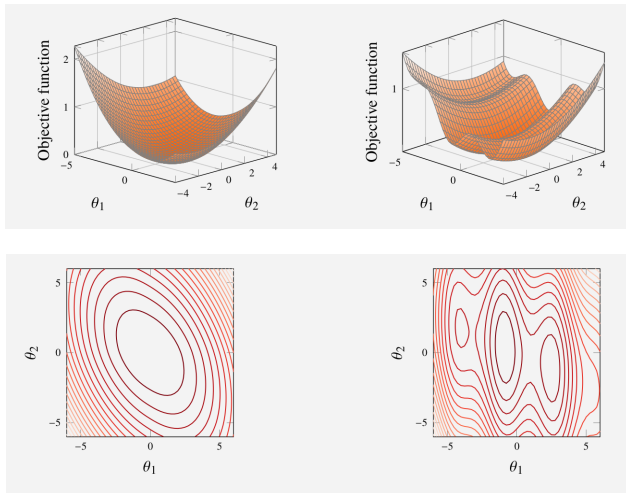
► Hyperparameter learning

$$\hat{\lambda} = \operatorname{argmin}_{\lambda} E_{\text{hold-out}}(\lambda)$$

■ Non-convexity and local minima.

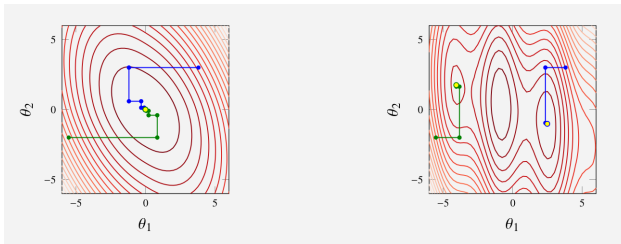


Parameter optimization



Coordinate descent

- Optimize $J(\theta_1, \theta_2, \dots, \theta_p)$ one coordinate θ_j at the time.
- L1 regularization for squared error loss: coordinate descent guaranteed to reach global minimum.



Gradient descent

- Idea: move in the direction of steepest descent, i.e. opposite direction of the **gradient**

$$\nabla_{\theta} J(\theta) = \left(\frac{\partial}{\partial \theta_1} J(\theta), \dots, \frac{\partial}{\partial \theta_p} J(\theta) \right)^{\top}$$

- Example logistic regression

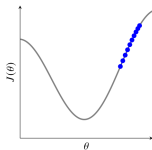
$$\nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{y_i x_i^{\top} \beta}} \right) y_i x_i$$

Algorithm 5.1: Gradient descent

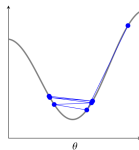
Input: Objective function $J(\theta)$, initial $\theta^{(0)}$, learning rate γ

Result: $\hat{\theta}$

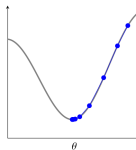
```
1 Set  $t \leftarrow 0$ 
2 while  $\|\theta^{(t)} - \theta^{(t-1)}\|$  not small enough do
3   | Update  $\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma \nabla_{\theta} J(\theta^{(t)})$ 
4   | Update  $t \leftarrow t + 1$ 
5 end
6 return  $\hat{\theta} \leftarrow \theta^{(t-1)}$ 
```



(a) Low learning rate $\gamma = 0.05$



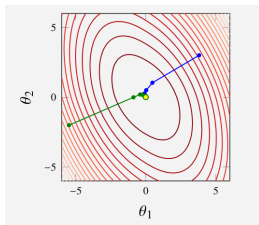
(b) High learning rate $\gamma = 1.2$



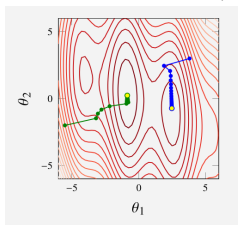
(c) Good learning rate $\gamma = 0.3$

Gradient descent

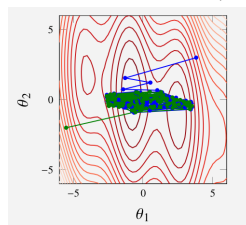
Convex



Nonconvex - good γ



Nonconvex - bad γ



Newton's method

- Gradient descent from first-order Taylor approximation of $J(\theta)$.
- **Newton's method** from second-order Taylor approx of $J(\theta)$:

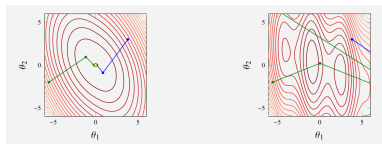
$$J(\theta + v) \approx \underbrace{J(\theta) + v^\top \nabla_\theta J(\theta) + \frac{1}{2} v^\top \nabla_\theta^2 J(\theta) v}_{s(\theta, v)}$$

- The approx $s(\theta, v)$ has minimum (solve $\nabla_v s(\theta, v) = 0$ for v)

$$v = -(\nabla_\theta^2 J(\theta))^{-1} \nabla_\theta J(\theta)$$

suggesting the iterative update

$$\theta^{(t+1)} = \theta^{(t)} - \left(\nabla_\theta^2 J(\theta^{(t)}) \right)^{-1} \nabla_\theta J(\theta^{(t)})$$



Newton's method with trust regions

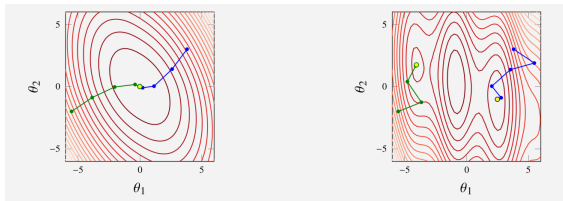
- Newton's method can be unstable. Hessian not always positive definite. Divergence.
- **Newton with trust regions**: don't allow large updates:

Algorithm 5.2: Trust-region Newton's method

Input: Objective function $J(\theta)$, initial $\theta^{(0)}$, trust region radius D

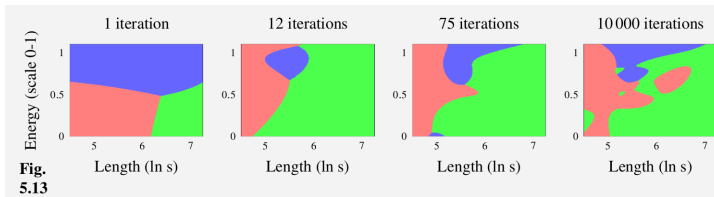
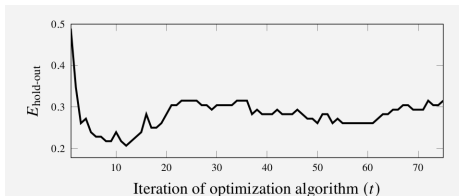
Result: $\hat{\theta}$

```
1 Set  $t \leftarrow 0$ 
2 while  $\|\theta^{(t)} - \theta^{(t-1)}\|$  not small enough do
3   Compute  $\mathbf{v} \leftarrow [\nabla_{\theta}^2 J(\theta^{(t)})]^{-1} [\nabla_{\theta} J(\theta^{(t)})]$ 
4   Compute  $\eta \leftarrow \frac{D}{\max(\|\mathbf{v}\|, D)}$ 
5   Update  $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{v}$ 
6   Update  $t \leftarrow t + 1$ 
7 end
8 return  $\hat{\theta} \leftarrow \theta^{(t-1)}$ 
```



Early stopping to avoid overfitting

- Song classification problem with multi-class logistic regression
- 20 degree polynomial with interactions. Overfitting!
- Monitor optimization progress on hold-out data.



Stochastic gradient descent

- **Large (independent) data**: gradient is a large sum. Costly!

$$\nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(x_i, y_i; \theta)$$

- **Stochastic gradient descent**:

- ▶ each iteration computes the gradient $\nabla_{\theta} J(\theta)$ on a random **mini-batch** of $n_b \ll n$ observations.
- ▶ each mini-batch gradient is an **unbiased estimator** of $\nabla_{\theta} J(\theta)$.
- ▶ mini-batch gradient descent. Learning rate γ_t in iteration t .

- **Epoch**: a complete sweep through all data. n/n_b iterations.

- Robbins-Monro **conditions for convergence**:

- ▶ $\sum_{t=0}^{\infty} \gamma_t = \infty$ [step length shouldn't vanish too fast]
- ▶ $\sum_{t=0}^{\infty} \gamma_t^2 < \infty$ [... but also not too slow]
- ▶ Fulfilled by $\gamma_t = \frac{1}{t^{\alpha}}$ for $\alpha \in (0.5, 1]$.

- However, often use lower cap: $\gamma_t \rightarrow \gamma_{\min} > 0$, e.g.

$$\gamma_t = \gamma_{\min} + (\gamma_{\max} - \gamma_{\min}) e^{-\frac{t}{\tau}}.$$

Stochastic gradient descent

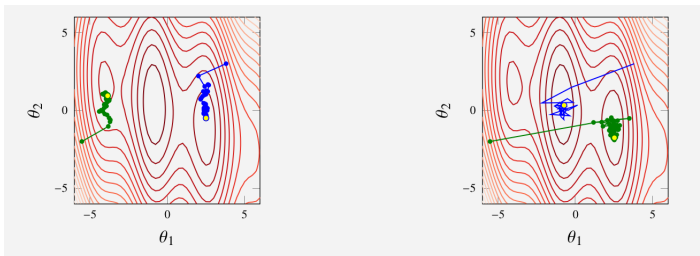
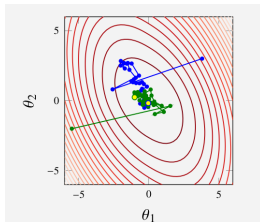
Algorithm 5.3: Stochastic gradient descent

Input: Objective function $J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, y_i, \theta)$, initial $\theta^{(0)}$, learning rate $\gamma^{(t)}$

Result: $\hat{\theta}$

```
1 Set  $t \leftarrow 0$ 
2 while Convergence criteria not met do
3   for  $i = 1, 2, \dots, E$  do
4     Randomly shuffle the training data  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ 
5     for  $j = 1, 2, \dots, \frac{n}{n_b}$  do
6       Approximate the gradient using the mini-batch  $\{(\mathbf{x}_i, y_i)\}_{i=(j-1)n_b+1}^{jn_b}$ ,
7        $\hat{\mathbf{d}}^{(t)} = \frac{1}{n_b} \sum_{i=(j-1)n_b+1}^{jn_b} \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta^{(t)})$ .
8       Update  $\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \hat{\mathbf{d}}^{(t)}$ 
9       Update  $t \leftarrow t + 1$ 
10    end
11  end
12 return  $\hat{\theta} \leftarrow \theta^{(t-1)}$ 
```

Stochastic gradient descent



Stochastic second order gradient methods

- Stochastic **Quasi-Newton** methods. Use past gradients to approximate Hessian info:

$$\text{learning rate:} \quad \gamma_t = \gamma(\nabla J_t, \nabla J_{t-1}, \dots, \nabla J_0)$$

$$\text{search directions:} \quad d_t = d(\nabla J_t, \nabla J_{t-1}, \dots, \nabla J_0)$$

- **Adaptive Stochastic Gradient Descent**

$$\theta^{(t+1)} = \theta^{(t)} - \gamma_t d_t$$

- The **ADAM** optimizer is based on exponential decay

$$d_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \nabla J_i$$

$$\gamma_t = \frac{\eta}{\sqrt{t}} \left((1 - \beta_2) \text{diag} \left(\sum_{i=1}^t \beta_2^{t-i} \|\nabla J_i\|^2 \right) \right)^{1/2}.$$

- β_1 and β_2 close to one is a common choice.

Hyperparameter learning

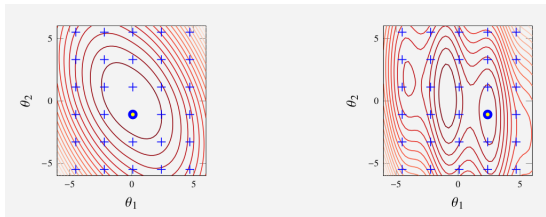
- Minimizing $E_{\text{hold-out}}(\lambda)$ by grid search.

Algorithm 5.4: Grid search for regularization parameter λ

Input: Training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, validation data $\{\mathbf{x}_j, y_j\}_{j=1}^{n_v}$

Result: $\hat{\lambda}$

```
1 for  $\lambda = 10^{-3}, 10^{-2}, \dots, 10^3$  (as an example) do
2   | Learn  $\hat{\theta}$  with regularization parameter  $\lambda$  from training data
3   | Compute error on validation data  $E_{\text{val}}(\lambda) \leftarrow \frac{1}{n_v} \sum_{j=1}^{n_v} (\hat{\mathbf{y}}(\mathbf{x}_j; \hat{\theta}) - y_j)^2$ 
4 end
5 return  $\hat{\lambda}$  as  $\arg \min_{\lambda} E_{\text{val}}(\lambda)$ 
```



- **Bayesian optimization** based on Gaussian processes (Lecture 8) when objective is costly and λ low-dim.