## Lec 17

Tuesday, November 12, 2019   10:59

# Recap: Neural Nets

Up until now, for linear models

$$\hat{f}(x) = g(\beta^T \varphi(x))$$

$\varphi$ given, $\beta$ learned

Big idea of NN: also learn $\varphi$

Vanila NN:

$$\overset{\mathbb{R}^P}{X} \longrightarrow \overset{\mathbb{R}^M}{Z = \varphi(x)} \longrightarrow \hat{f}(x) = \beta^T Z$$

$$\varphi_m(x) = \sigma(\alpha_{0m} + \alpha_m^T x) \quad m = 1, \cdots, M$$

Given $\alpha$'s, $\beta$'s, the steps to compute the output of the vanila neural network:

— Get input $x \in \mathbb{R}^P$

— Compute hidden layer $z_j = \sigma(\alpha_{j0} + \alpha_j^T x)$

— Compute outputs $T_k = \beta_{k0} + \beta_k^T z$

   — For regression, just return outputs

   — For classification, apply softmax & return probabilities

$$\hat{P}_k = \frac{e^{T_k}}{\sum_{k'} e^{T_{k'}}}$$

The params of this model are $\alpha$'s, $\beta$'s

Want to find params that min loss on training data / max likelihood

For regression: min sum of squared errors

$$\min_{\alpha, \beta} \sum_i \sum_k (Y_{ik} - T_{ik})^2 \underbrace{\quad\quad}_{\text{dependence on } x_i, \alpha, \beta}$$

For classification: max likelihood / min cross entropy

$$\min_{\alpha, \beta} -\sum_i \sum_k Y_{ik} \log \hat{P}_{ik} \underbrace{\quad\quad}_{\text{depends on } \alpha, \beta, x_i}$$

Can add ridge regularization:

$$\min_{\text{params}} \text{loss} + \lambda \underbrace{\| \text{params} \|_2^2}_{}$$

$$\alpha_{11}^2 + \alpha_{12}^2 + \cdots + \alpha_{1p}^2$$
$$+ \alpha_{21}^2 + \cdots + \beta_{11}^2 + \cdots$$

called "weight decay"

How to actually solve $\min_{\text{params}} \text{loss}$ ?

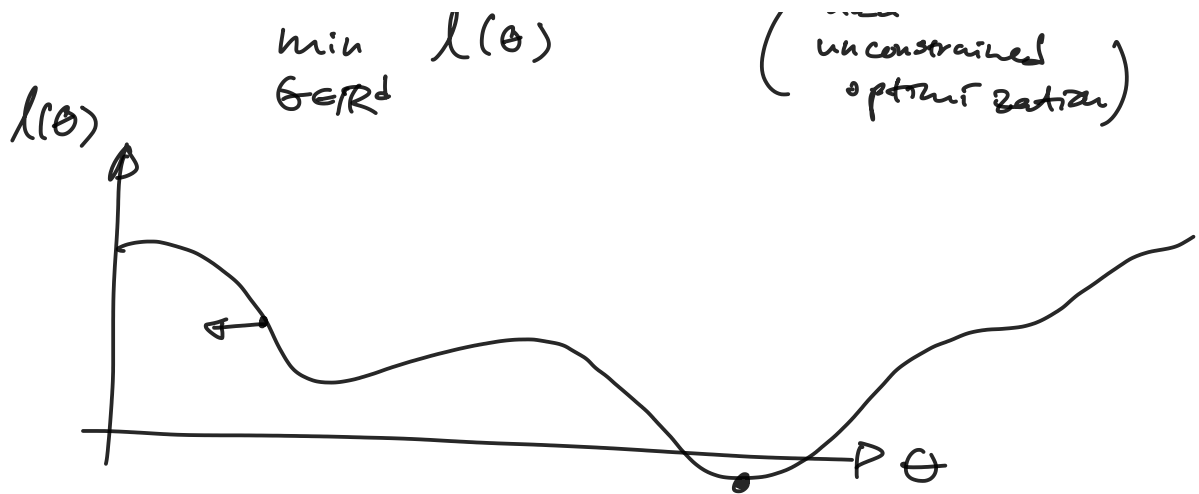Put NN to the side.

Let's tackle general problem

Gradient Descent & Variants

Given a function $\ell(\theta), \quad \theta \in \mathbb{R}^d$

$$\ell : \mathbb{R}^d \longrightarrow \mathbb{R}$$

Want to solve

$$\min_{\theta \in \mathbb{R}^d} \ell(\theta)$$

( unconstrained optimization )
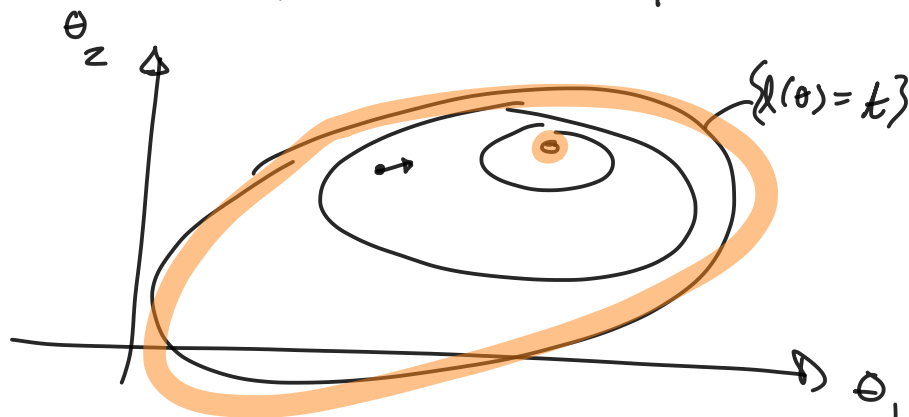


Recall: $\nabla \ell(\theta) = \begin{pmatrix} \partial \ell(\theta)/\partial \theta_1 \\ \vdots \\ \partial \ell(\theta)/\partial \theta_d \end{pmatrix} \in \mathbb{R}^d$

$$\nabla \ell : \mathbb{R}^d \longrightarrow \mathbb{R}^d$$



$\{\ell(\theta) = t\}$

Gradient = direction of highest ascent

To min, want to go down

− Grad = direction of greatest descent

Idea: to min a fn (i.e. find it's lowest pt in the graph), walk in the direction of greatest descent

## Vanilla Gradient Descent Algo

Inputs: − Starting pt $\theta_0 \in \mathbb{R}^d$

         − Gradient fn $\nabla \ell(\theta)$

– Step size $\alpha$

— Tolerance $\tau, m$

For $t = 1, 2, \cdots$

$$\theta_t \leftarrow \theta_{t-1} - \alpha \nabla \ell(\theta_{t-1})$$

If $\|\theta_{t-i} - \theta_{t-i-1}\| < \tau \qquad \forall i = 0, \cdots, m-1$

Return $\theta_t$

Variations on choosing stepsize $\alpha$

One variation: Line search

$$\alpha_t = \underset{\alpha_t \geq 0}{\text{argmin}} \; \ell\left(\theta_{t-1} - \alpha \nabla \ell(\theta_{t-1})\right)$$

$$\theta_t = \theta_{t-1} - \alpha_t \nabla \ell(\theta_{t-1})$$

Now suppose $\ell(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell_i(\theta)$

e.g. $\ell_i(\theta) = \left(Y_i - \hat{F}(x_i; \theta)\right)^2$

$$\nabla \ell(\theta) = \frac{1}{n} \nabla \sum_i \ell_i(\theta)$$

$$= \frac{1}{n} \sum_i \nabla \ell_i(\theta)$$

In words: "the gradient of the average
is the average of the gradients"

To run gradient descent we'd need
to compute $\nabla \ell_i(\theta)$ for $i = 1, \cdots, n$
in order to get $\nabla \ell(\theta)$

~~In order to~~ get $\nabla \lambda(\theta)$

Just for one step

Maybe too much when $n$ is very large

## Stochastic Gradient Descent

1. Draw $i \sim unif\{1, \cdots, n\}$
2. $\theta_t \leftarrow \theta_{t-1} - \alpha \nabla l_i(\theta_{t-1})$

Other var.