

# Why Nonlinear Optimization?

Supervised Learning

$$x_i \in \mathbb{R}^{N_x}$$

$(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$  — training data  
input      output labels  
 $(x, y)$

$$\sim p(x, y)$$

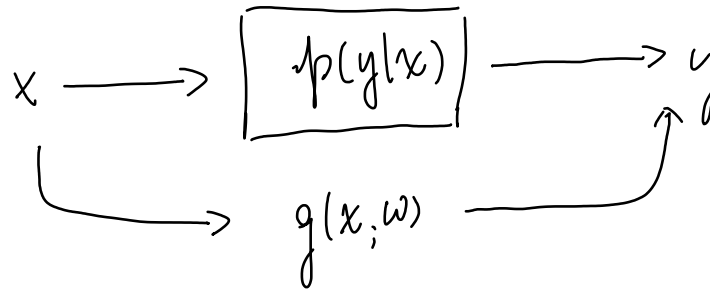
$$x \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} \left. \vphantom{\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}} \right\} \rightarrow y = 1$$

$(x, y)$        $p(x, y)$

Recognition

$$p(y|x) = \frac{p(x, y)}{\int p(x, y) dy} = p_x(x)$$

training data  $(x_1, y_1) \dots (x_n, y_n)$  i.i.d.  $\sim p(x, y)$



$$g(x; w) \approx y$$
$$\text{error} = |g(x; w) - y|$$

$$\mathbb{E}_{(x,y) \sim p(x,y)} |g(x;w) - y| = \int_Y \int_X |g(x;w) - y| p(x,y) dx dy \quad \text{"population loss"}$$

$$= L(w)$$

Learning  $(w^*) = \min_w L(w)$

Framework

$$(x, y) \sim p(x, y)$$

$$g(x; w) = wx$$

$$w^T x = w_1 x_1 + \dots + w_n x_n$$

$$L(w) = \mathbb{E}_{(x,y) \sim p(x,y)} l(g(x,w); y)$$

$l(g; y)$  — loss function

$$x \longrightarrow g(x; \omega^*) - y \neq 0$$

$$\underline{g(x; \omega)}$$

$$\mathbb{E}_{(x,y) \sim p(x,y)} \ell(g(x; \omega), y) \approx \frac{\ell(g(x_1; \omega), y_1) + \dots + \ell(g(x_n; \omega), y_n)}{\textcircled{n}}$$

"empirical loss"

$$L_n(\omega) = \frac{1}{n} \sum_{i=1}^n \ell(g(x_i; \omega), y_i) \xrightarrow{n \rightarrow \infty} \mathbb{E}_{\substack{(x,y) \\ \sim p(x,y)}} \ell(g(x; \omega), y)$$

$(x_1, y_1) \dots (x_n, y_n)$  training data

$$\arg \min_w \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, w), y_i) = w_n^* \quad \text{"Optimization"}$$

Empirical Risk Minimization (ERM).

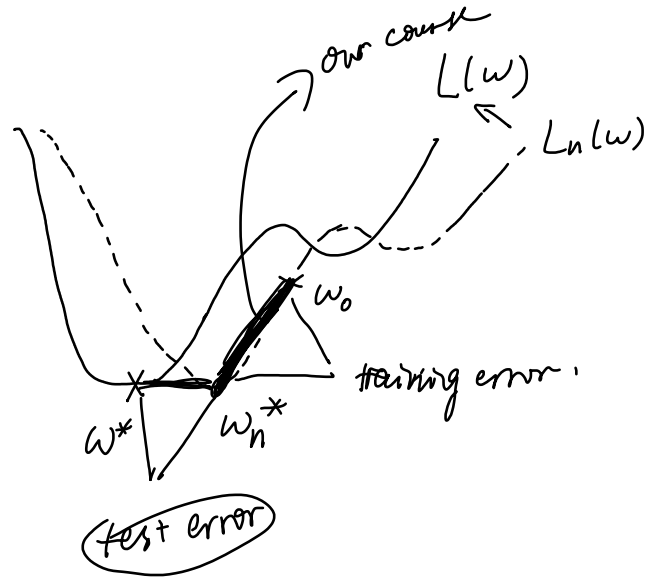
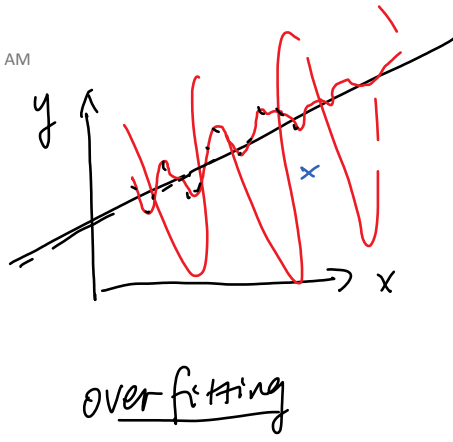
$$w_n^* \xrightarrow{n \rightarrow \infty} w^*$$

$$\arg \min L_n(w) \quad \arg \min \underline{L(w)}$$

$$L_n \longrightarrow L$$

$$\mathbb{E}_{(x,y)} \ell(f(x, w_n^*); y)$$

"Statistics"  $(x, y) \sim p(x, y)$   
test data



Empirical Risk Minimization  $L_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(g(x_i; w), y_i)$

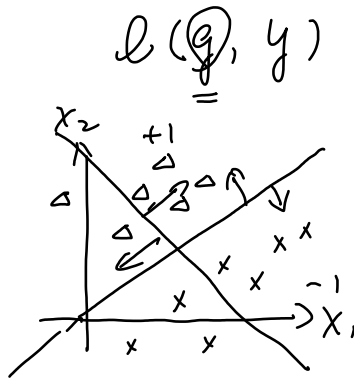
Learning Problem  $\min_w L_n(w) = w_n^*$

Concrete Stuff

"Loss function"

$x = (x_1, x_2)$

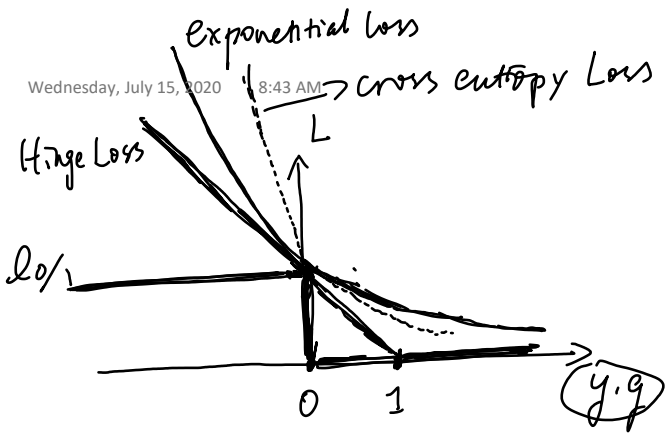
$y = (1, -1)$



classification  
estimation

0-1 loss function

$$\ell_{0/1}(y, g) = \begin{cases} 1 & \text{if } y \cdot g < 0 \\ 0 & \text{otherwise} \end{cases}$$



$$L = |d| = \begin{cases} 0 & \text{if } y \cdot g > 0 \\ 1 & \text{if } y \cdot g \leq 0 \end{cases}$$

(1) Hinge Loss

$$L(g, y) = \max(0, 1 - yg)$$

(2) Exponential Loss

$$L(g, y) = e^{-y \cdot g}$$

(3) Cross Entropy Loss

$$L(g, y) = - \left( \mathbb{1}_{y=1} \ln \left( \frac{e^g}{e^g + e^{-g}} \right) + \mathbb{1}_{y=-1} \ln \left( \frac{e^{-g}}{e^g + e^{-g}} \right) \right) = - \ln \left( \frac{e^{y \cdot g}}{e^{y \cdot g} + e^{-y \cdot g}} \right)$$



$z_1 \dots z_n$

$$\frac{e^{z_j}}{e^{z_1} + \dots + e^{z_n}}$$

$$y = (1, 2, 3, 4)$$

$$g = (1, 3, 1, 4)$$

$y - g$

Estimation (Regression)

$$L(y, g) = |y - g|_2^2 \quad (L^2 - \text{Loss})$$

$$|x|_2^2 = x_1^2 + \dots + x_d^2$$

$$L(y, g) = |y - g|_1 \quad (L^1 - \text{Loss})$$

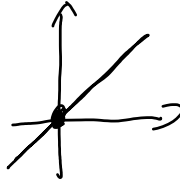
$$|x|_1 = |x_1| + \dots + |x_d|$$

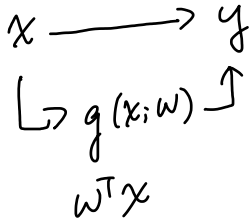
$$L(y, g) = |y - g|_0 \quad (L^0 - \text{Loss}) \quad |x|_0 = \#\{i : x_i \neq 0\}$$

$$x = (x_1, \dots, x_d) \quad x = (x_1, \dots, x_d)$$

# Learning Functions

## (1) Linear Regression

$$g(x; w) = \frac{w^T x}{1} = w_1 x_1 + \dots + w_d x_d$$




$$w^T x + \beta$$

$$\begin{cases} \tilde{x} = (x, 1) \\ \tilde{w} = (w, \beta) \end{cases}$$

$$x = (x_1 \dots x_d)$$

$$w = (w_1 \dots w_d)$$

$$w^T x + \beta = \tilde{w}^T \tilde{x}$$

# Least Squares Problem

$$(x_j, y_j) \quad \frac{x_j \in \mathbb{R}^d}{y_j \in \mathbb{R}} \quad y = (y_1, \dots, y_n)$$

$$L(g, y) = \frac{1}{2} (g - y)^2$$

$$g(x; w) = w^T x \quad w \in \mathbb{R}^d \quad w^T x = x^T w$$

$$L(g(x_j; w), y_j) = \frac{1}{2} (w^T x_j - y_j)^2$$

$$\text{ERM} \quad w_n^* = \min_w \frac{1}{n} \sum_{j=1}^n \frac{1}{2} \underline{\underline{(w^T x_j - y_j)^2}}$$

$$A = (x_1^T \dots x_n^T) \quad (x_1^T w - y_1, \dots, x_n^T w - y_n)$$

$$= Aw - y$$

ERM  $w_n^* = \arg \min_w \frac{1}{2n} \|Aw - y\|_2^2$

Big data "High-dimensional Stat."  $d \gg 1$

$(w = |w_1| + \dots + |w_d|)$

$L^1$ -regularization

(LASSO Least Absolute Shrinkage and Selection Operator)

$$\min_w \left[ \frac{1}{2n} \|Aw - y\|_2^2 + \lambda \sum |w_i| \right]$$

$\lambda > 0$

tuning parameter

penalization

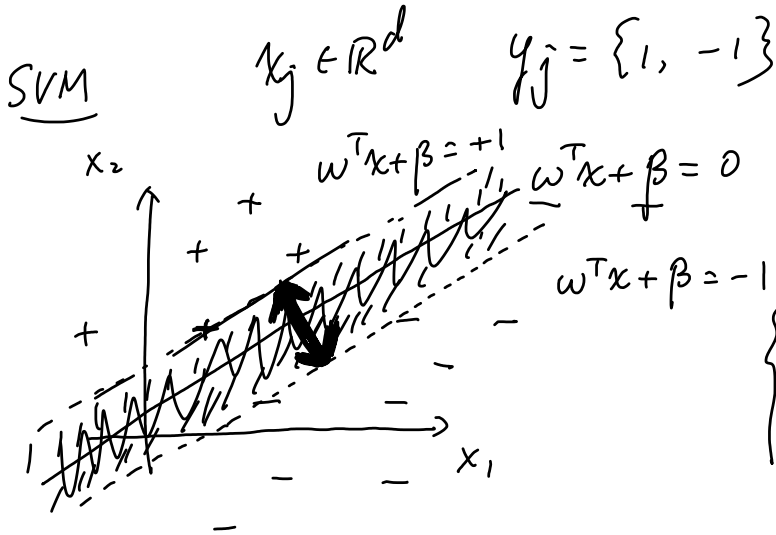
Regularized ERM

$$w_n^* = \arg \min_w \frac{1}{n} \sum_{i=1}^n \left[ \underbrace{L(g(x_i, w), y_i)}_{l_i(w)} + \lambda R(w) \right]$$

$$L_n(w) = \frac{1}{n} \sum_{i=1}^n l_i(w)$$

# SVM — Support Vector Machines

Neural Network



"hyperplane"

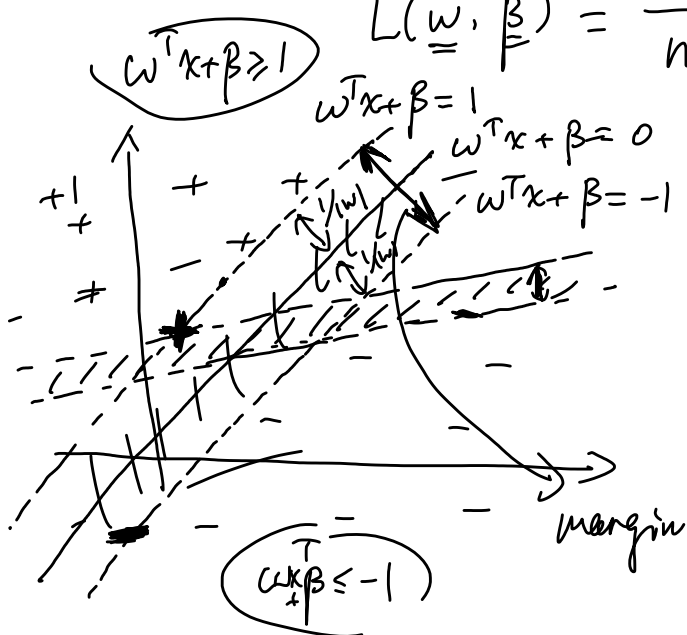
Seek for  $w \in \mathbb{R}^d$   $\beta \in \mathbb{R}$  s.t.

$$\begin{cases} w^T x_j + \beta \geq 1 & \text{if } y_j = 1 \\ w^T x_j + \beta \leq -1 & \text{if } y_j = -1 \end{cases}$$

$$L(y, \hat{y}) = \max(0, 1 - y\hat{y})$$

Hinge Loss

$$L(\underline{w}, \underline{\beta}) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_j (x_j^T \underline{w} + \beta))$$



$$x \in \mathbb{R}^d \quad \Pi: w^T x + \beta = 0$$

$$\text{dist}(x, \Pi) = \frac{|w^T x + \beta|}{|w|}$$

$$\text{width of the strip} = d = \frac{2}{|w|}$$

SVM

$$\max_{w \in \mathbb{R}^d, \beta \in \mathbb{R}} \frac{2}{|w|}$$

s.t.  $y_j (w^T x_j + \beta) \geq 1$   
for  $j = 1 \dots n$

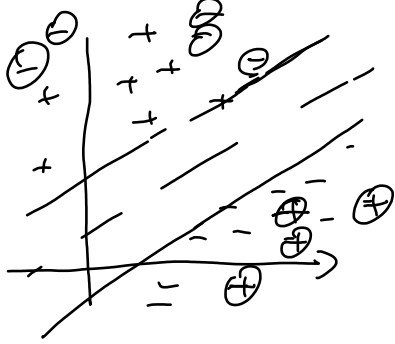
$$\min_{w \in \mathbb{R}^d, \beta \in \mathbb{R}} \frac{1}{2} |w|^2$$

s.t.  $y_j (w^T x_j + \beta) \geq 1$   
for  $j = 1 \dots n$



⊖

Wednesday, July 15, 2020 9:36 AM



"Soft margin"

$$\min_{\substack{w \in \mathbb{R}^d \\ \beta \in \mathbb{R}}} \left[ \frac{1}{2} \|w\|^2 + C \sum_{j=1}^n \max(0, 1 - (y_j (w^T x_j + \beta))) \right]$$

Hinge

$$\max_{\substack{w \in \mathbb{R}^d \\ \beta \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + C \sum_{j=1}^n \xi_j \quad \text{s.t.} \quad \underline{y_j (w^T x_j + \beta)} \geq 1 - \xi_j$$

↑  
slack variables

$\xi_j \geq 0$

# Neural Networks

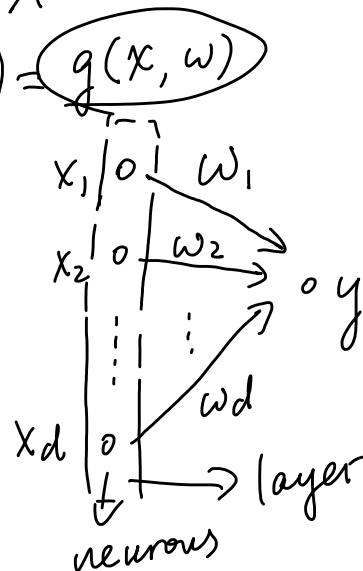
Least Squares

$$y \approx \underbrace{A}_{(x)} x \iff \min_A \frac{1}{2} |y - Ax|^2$$

$$y \approx \underline{w_1 x_1 + w_2 x_2 + \dots + w_d x_d} = g(x, w)$$

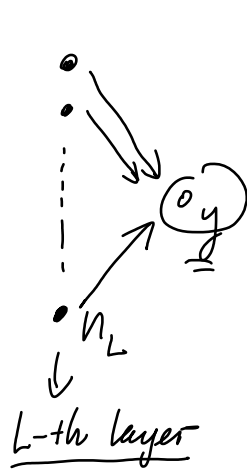
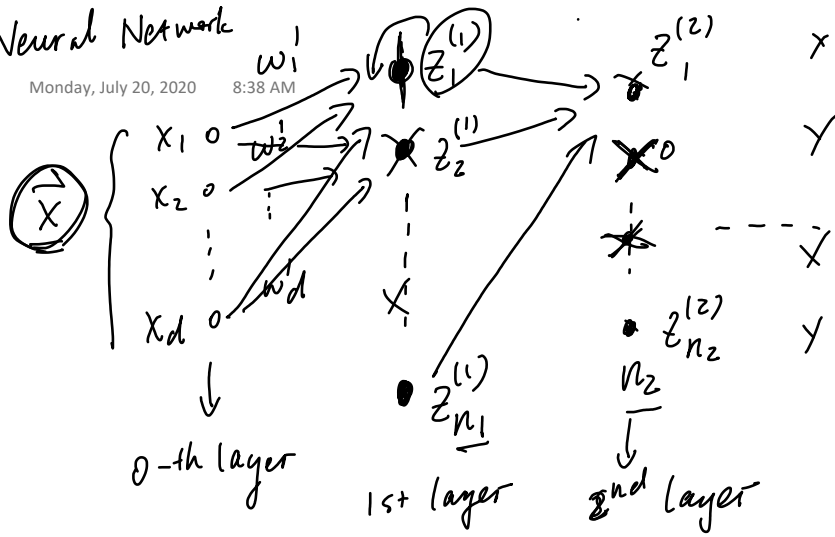
$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

$$A = (w_1, \dots, w_d)$$



# Neural Network

Monday, July 20, 2020 8:38 AM



$$z^{(1)} = W^{(1)}x + b^{(1)}$$

$$z^{(2)} = W^{(2)}\sigma(z^{(1)}) + b^{(2)}$$

$$\vdots$$

$$z^{(L)} = W^{(L)}\sigma(z^{(L-1)}) + b^{(L)}$$

$$y = a^T z^{(L)}$$

$$= \langle a, z^{(L)} \rangle$$

$$z_k^{(1)} = w_1^k x_1 + w_2^k x_2 + \dots + w_d^k x_d \quad k=1, \dots, n_1$$

$$z^{(1)} = \begin{pmatrix} z_1^{(1)} \\ z_2^{(1)} \\ \vdots \\ z_{n_1}^{(1)} \end{pmatrix} = \begin{pmatrix} w_1 & \dots & w_d \\ w_1^k & \dots & w_d^k \\ \vdots & \dots & \vdots \\ w_{n_1} & \dots & w_d \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \rightarrow W^{(1)} \text{ --- weight.}$$

$$b^{(1)} \in \mathbb{R}^{n_1} \text{ --- bias}$$

$$y = a^T \left[ W^{(L)} \left( \dots \sigma \left[ W^{(3)} \left( W^{(2)} \left[ \sigma \left( W^{(1)} x + b^{(1)} \right) \right] + b^{(2)} \right) \right] + b^{(3)} \right) \dots \right] + b^{(L)}$$

$$= a^T \underbrace{W^{(L)} W^{(L-1)} \dots W^{(1)}}_W x + \dots = W x + b$$

$\min_{W^{(i)}}$

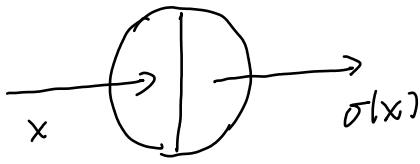
$$\frac{1}{2} \left| y - a^T \underbrace{\hat{W}^{(L)} \hat{W}^{(L-1)} \dots \hat{W}^{(1)}}_W x \right|^2$$

↓  
adaptive

"matrix factorization"

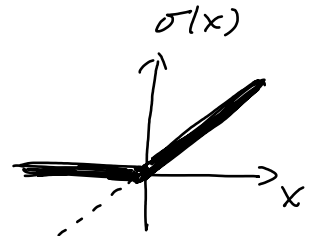
$$x \xrightarrow{W} y$$

$$\begin{bmatrix} 0 & 0 & \dots \\ & \mathcal{R} & \\ 0 & 0 & \end{bmatrix} \approx \begin{bmatrix} | \\ | \\ | \end{bmatrix} =$$



$$\sigma(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

↓  
activation function



ReLU  
(Rectified Linear Unit)

$$\sigma(x_1, \dots, x_d) = (\sigma(x_1), \dots, \sigma(x_d))$$

tanh

Neural-Network

"Linear" + "Activation"

Sigmoid

I. Goodfellow --- "Deep Learning" MIT

$$y \approx a^T \left( \underbrace{\sigma(\underline{W}^{(L)}) \sigma(\underline{W}^{(L-1)}) \dots \sigma(\underline{W}^{(1)} x + b^{(1)}) \dots + b^{(L-1)}}_{g(x, w)} + b^{(L)} \right)$$

$g(x, w)$

$$\omega = \begin{pmatrix} W^{(1)} & \dots & W^{(L)} \\ b^{(1)} & \dots & b^{(L)} \end{pmatrix}$$

$$\min_w \frac{1}{2} (y - g(x; \underline{w}))^2$$

"Interpretable Machine/Deep Learning"

ReLU  $\sigma(x_1, \dots, x_d) = (x_1, x_2, 0, x_3, \dots, 0, \dots, x_d)$  } Learning  
 "adaptive piecewise linear model"

$$\underline{\underline{F(w)}} = \frac{1}{2} (y - g(x; \underline{w}))^2 \quad \text{"loss landscape"}$$

$w_1, w_2$

