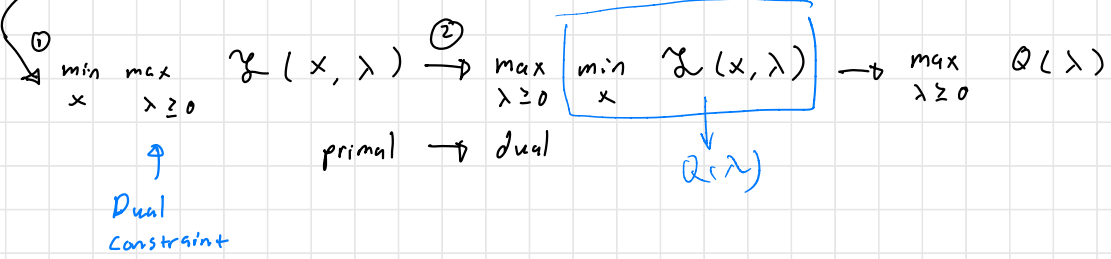# Lecture 15

- I.O.V: Constrained Optimization

$$\min \quad f(x)$$
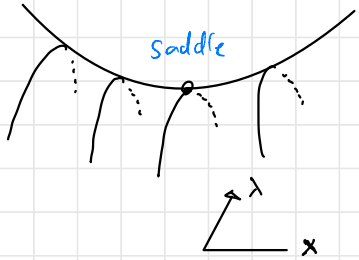$$\text{s.t. } g(x) \leq 0 \quad \longleftarrow \text{ primal constraint}$$

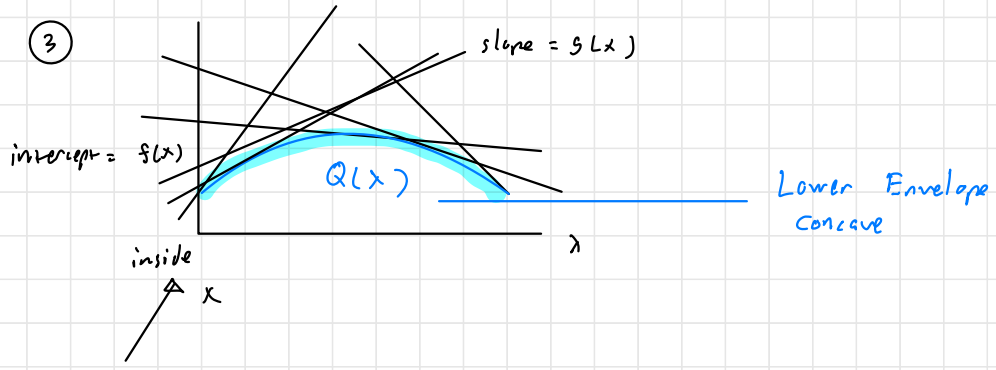$$\mathcal{L}(x, \lambda) = f(x) + \lambda g(x)$$

① $\underset{x}{\min} \underset{\lambda \geq 0}{\max} \mathcal{L}(x, \lambda) \xrightarrow{②} \underset{\lambda \geq 0}{\max} \boxed{\underset{x}{\min} \mathcal{L}(x, \lambda)} \longrightarrow \underset{\lambda \geq 0}{\max} \quad Q(\lambda)$

$\uparrow$ Dual Constraint

primal $\rightarrow$ dual

$Q(\lambda)$

Note: ①

$$\underset{\lambda \geq 0}{\max} \quad \lambda \boxed{g(x)} \longrightarrow \infty$$
$$> 0$$

② Von Neumann

convex - concave



Saddle

$\lambda$

$x$

The Lagrangian is linear in $\lambda$

③



slope = $g(x)$

intercept = $f(x)$

$Q(\lambda)$

$\lambda$

inside

$x$

Lower Envelope Concave

$f(x)$ — slope $g(x) = 0$

$Q(\lambda)$

$\lambda > 0$

0   $\lambda$

$\lambda = 0$

$f(x)$  slope $g(x) < 0$

$Q(\lambda)$

0   $\lambda$

○ What if we had

$$\min \ f(x)$$
$$\text{s.t.} \ g(x) = 0$$

Then no need for $\lambda \gtreqless 0$

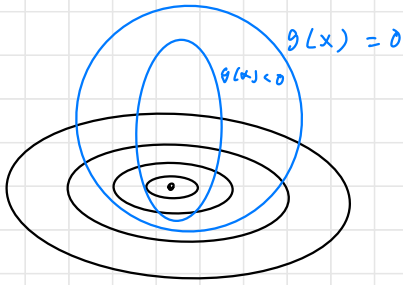$$\min_x \max_\lambda \ \mathcal{L}(x, \lambda)$$

$$\max_\lambda \ \lambda g(x)$$

$\lambda \to \infty$

if $g(x) > 0$ → $+\infty$

$\lambda \to -\infty$

if $g(x) < 0$ → $+\infty$

③

$g(x) = 0$

$\lambda$

- Consider the Contour Plot:



$g(x) = 0$

$g(x) < 0$

$g(x) = 0$

$g(x) < 0$

$\lambda = 0$

$f'(x) = -\lambda g'(x)$

$f'(x) + \lambda g'(x) = 0$

$$\frac{d}{dx} \mathcal{L}(x, \lambda) = 0$$

○ KKT conditions

(1) primal constraints
(2) Dual constraints
(3) Stationarity $\frac{d}{dx} \mathcal{L} = 0$
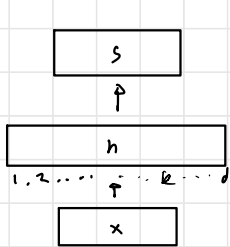(4) Complementary slackness ; $\lambda g(x) = 0$

- Deep Learning :

  -

    - ReLU :     $f(x)$  piecewise  linear
    - Sigmoid :   $f(x)$   piecewise  constant

| S |
|---|

$\uparrow$
p

| h |
|---|

1, 2, ... , $\uparrow$ , k ... d

| x |
|---|

$$S = f(x) = \frac{1}{\sqrt{d}} \sum_{k=1}^{d} \beta_k \, \sigma(x^T \alpha_k)$$

- Initialization :     $\begin{pmatrix} \alpha_k \\ \beta_k \end{pmatrix} \overset{iid}{\sim} P_0(\alpha_k, \beta_k)$  ,  $k = 1, \dots, d$



$x \quad x'$

$k(x, x') \sim cov(f(x), f(x')) = cov\left(\frac{1}{\sqrt{d}} \sum \beta_k \sigma(x^T \alpha_k), \frac{1}{\sqrt{d}} \sum \beta_k \sigma(x'^T \alpha_k)\right)$

$\quad = \frac{1}{d} \sum_{k=1}^{d} E\left(\beta_k^2 \, \sigma(x^T \alpha_k) \, \sigma(x'^T \alpha_k)\right)$

$\quad = E\left(\beta_k^2 \, \sigma(x^T \alpha_k) \, \sigma(x'^T \alpha_k)\right)$     Central limit Theorem

So  $f(x) \sim NNGP(\quad) \xrightarrow{d \to \infty}$  CLT converges to GP

Neural net Gaussian process

1. Gradient Descent Learning

  - freeze $\alpha_k \overset{iid}{\sim} P_0(\alpha)$

    $$f(x) = h(x)^T \beta$$

    $$h_k = \frac{1}{\sqrt{d}} \sigma(x^T \alpha_k) \quad \text{random feature}$$

    learn $\beta$

$$\text{Loss}(\beta) = \frac{1}{2n} \sum_{i=1}^{n} (x_i - h(x_i)^T \beta)^2$$

  $\beta_0 = 0$

  $\beta_{v+1} \quad > \quad \beta_b + \eta_t \frac{1}{n} \sum_{i=1}^{n} \underbrace{(x_i - h(x_i T) \beta_b)}_{e_i} h(x_i)$

  $\longrightarrow \hat{\beta} = \sum_{i=1}^{n} c_i h(x_i) \quad \text{representer.}$

  $\hat{f}(x) = h(x)^T \hat{\beta} = \sum_{i=1}^{n} c_i k(x_i, x) \quad$ <span style="color:red">kernel regression</span>
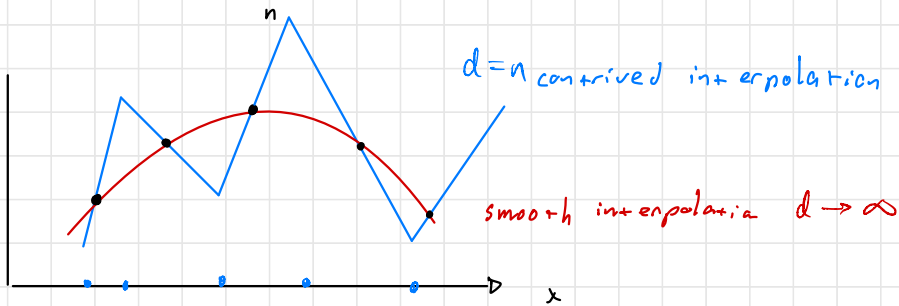  $\downarrow$ $\hspace{6cm}$ <span style="color:red">$\lambda \longrightarrow 0$</span>
  <span style="color:blue">query</span>

• Over parameterization : $\quad d \longrightarrow \infty$

- Double Descent

Error



- Consider  unfreeze  $\alpha_k \overset{iid}{\sim} p_0(\alpha)$

$$\theta = [\theta_k = (\alpha_k, \beta_k) , k = 1, \ldots, d)$$

initialization  $\theta_k^0 = (\alpha_k^0, \beta_k^0) = p_0(\alpha, \beta)$

$$f_\theta(x) = f_{\theta_0}(x) + \langle \theta - \theta_0, \nabla_\theta f_{\theta_0}(x)\rangle + \quad 2^{nd} \text{ order}$$

$\downarrow$
query
fixed

$$f_\theta(x) = \frac{1}{\sqrt{d}} \sum_{k=1}^{d} \beta_k \, \sigma(x^\top d_k) = \frac{1}{\sqrt{d}} \sum_{k=1}^{d} h_{\theta_k}(x)$$

$$= f_{\theta_0}(x) + \frac{1}{\sqrt{d}} \sum_{k=1}^{d} \nabla_{\theta_k} h_{\theta_k^0}(x)(\theta_k - \theta_k^0) + \frac{1}{2} \frac{1}{\sqrt{d}} \sum_{k=1}^{d} \nabla_{\theta_k}^2 h_{\theta_k^0}(x)(\theta_k - \theta_k^0)^2$$

$$\downarrow$$
$$0$$

$$\Delta = f_\theta(x) - f_{\theta_0}(x) \quad \text{so} \quad (\theta_k - \theta_k^0) \propto \frac{1}{\sqrt{d}} \quad \textcolor{red}{\text{perturbation}}$$

$$\text{also} \quad (\theta_k - \theta_k^0)^2 \propto \frac{1}{d}$$

$$\underline{\text{Thus the second term}} \longrightarrow 0$$

○ Note : As $d \to \infty$ The model behaves as a linear model

$$\textcolor{blue}{(\text{Tangent Model})}$$

$$f_\theta(x) = f_{\theta_0}(x) + \langle \theta - \theta_0, \nabla_\theta f_{\theta_0}(x) \rangle + 2^{nd} \text{ order}$$

$$\underbrace{\phantom{\langle \theta - \theta_0, \nabla_\theta f_{\theta_0}(x) \rangle}}$$

$$\textcolor{red}{\tilde{h}(x)}$$

Neural Tangent Kernel Regime

$$\langle \tilde{h}(x), \tilde{h}(x') \rangle = \frac{1}{d} \sum_{k=1}^{d} \nabla_{\theta_k} h_{\theta_k^0}(x) \cdot \nabla_{\theta_k} h_{\theta_k^0}(x')$$
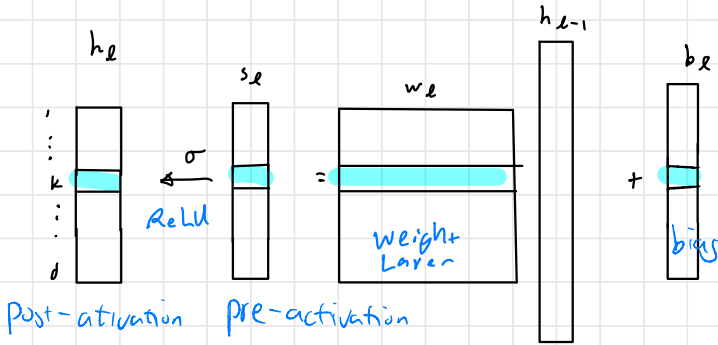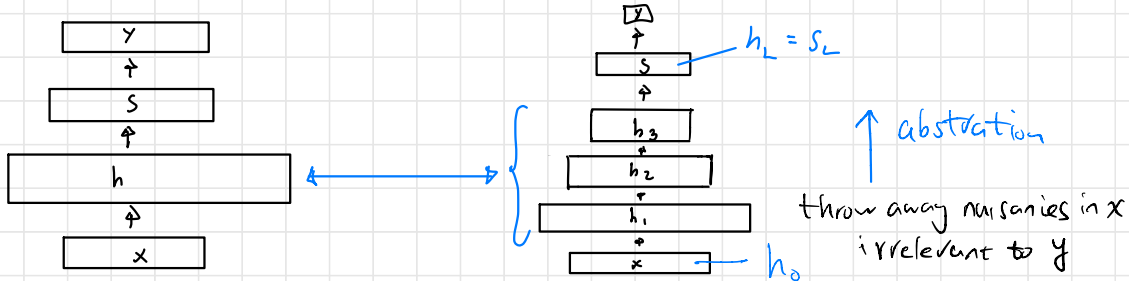
$$\longrightarrow E_{\theta_k} \left( \phantom{xxxxxx} \right)$$

Law of large #

- Multi-Layer Perceptron (MLP)

Recall:

MLP:



$h_L = s_L$

↑ abstration

throw away nuisances in $x$
irrelevant to $y$

$h_0$



Post-activation    Pre-activation    weight Layer    bias

$$h_{l,k} = \sigma(s_{l,k})$$

$$s_l = w_l h_{l-1} + b_l \Rightarrow s_{l,k} = w_{l,k} h_{l-1} + b_{l,k}$$

$$\theta = (w_l, b_l, \quad l = 1, \dots, L)$$

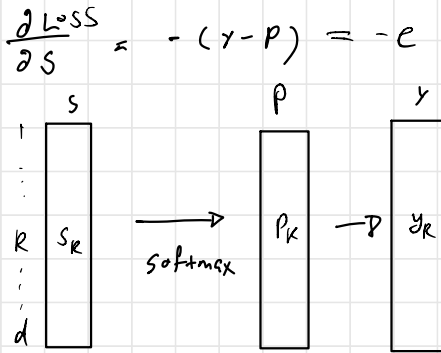- Regression: $\text{Loss} = \frac{1}{2} |Y - s|^2$

$$\frac{\partial \text{Loss}}{\partial s} = -(Y - s) = -e$$

- Classification:

$$P(Y|s) = \frac{e^{\langle Y, s \rangle}}{Z} = \frac{e^{Y_k s_u}}{\sum_{c=1} e^{Y_c s_c}}$$

$$\text{Loss} = -\log P(Y|s) \quad \text{cross-entropy}$$

$$\frac{\partial \text{Loss}}{\partial s} = -(Y - P) = -e$$



$s \xrightarrow{\text{softmax}} P \to y$

generalized linear model
(GLM)

- Error back-prop

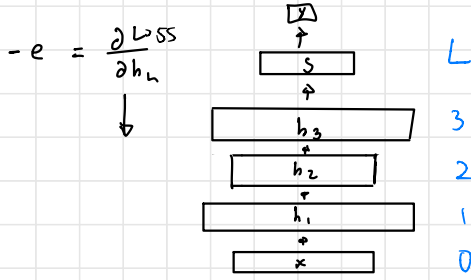$$-e = \frac{\partial \text{Loss}}{\partial h_n}$$

Loss



L

3

2

1

0

$$\frac{\partial \text{Loss}}{\partial h_\ell} \longrightarrow \frac{\partial \text{Loss}}{\partial w_\ell}, \frac{\partial \text{Loss}}{\partial b_\ell}$$

$$\frac{\partial \text{Loss}}{\partial h_{\ell-1}}$$

$$\frac{\partial \text{Loss}}{\partial h_{\ell-1}}^T = \sum_{k=1}^{d} \frac{\partial \text{Loss}}{\partial h_{\ell,k}} \cdot \frac{\partial h_{\ell,k}}{\partial s_{\ell,k}} \frac{\partial s_{\ell,k}}{\partial h_{\ell-1}}^T$$

$$= \sum_{k=1}^{d} \frac{\partial \text{Loss}}{\partial h_{\ell,k}} \sigma'_k(s_{\ell,k}) w_{\ell,k}$$

$$= \left( \cdots \quad \frac{\partial \text{Loss}}{\partial h_{\ell,k}} \quad \cdots \right) \begin{pmatrix} \vdots \\ \sigma'_k(s_{\ell,k}) w_{\ell,k} \\ \vdots \end{pmatrix} \begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ d \end{matrix}$$

$$\begin{matrix} 1 & & k & \cdots & d \end{matrix}$$

$$= \frac{\partial \text{Loss}}{\partial h_\ell}^T \quad \sigma'_\ell \odot w_\ell$$

row-wise multiply

$$\sigma'_\ell = \begin{pmatrix} \vdots \\ \sigma'_k(s_{\ell,k}) \\ \vdots \end{pmatrix} \begin{matrix} 1 \\ \vdots \\ k \\ \vdots \\ d \end{matrix}$$

- $$\frac{\partial \text{Loss}}{\partial w_{\ell,k}} = \frac{\partial \text{Loss}}{\partial h_{\ell,k}} \frac{\partial h_{\ell,k}}{\partial s_{\ell,k}} \cdot \frac{\partial s_{\ell,k}}{\partial w_{\ell,k}}$$

$$= \frac{\partial \text{Loss}}{\partial h_{\ell,k}} \sigma'(s_{\ell,k}) h_{\ell-1}^T$$

$$= \frac{\partial Loss}{\partial w} = \begin{pmatrix} \frac{\partial Loss}{\partial w_{\ell, k}} \\ \vdots \\ \cdots \end{pmatrix}_{k}^{\iota} = \begin{pmatrix} \sigma'(s_{\ell, k}) \frac{\partial Loss}{\partial h_{k, \ell}} \end{pmatrix} h_{\ell-1}^{T} =$$

$$= \sigma_{\ell}' \odot \frac{\partial Loss}{\partial h_{\ell}} \cdot h_{\ell-1}^{T}$$