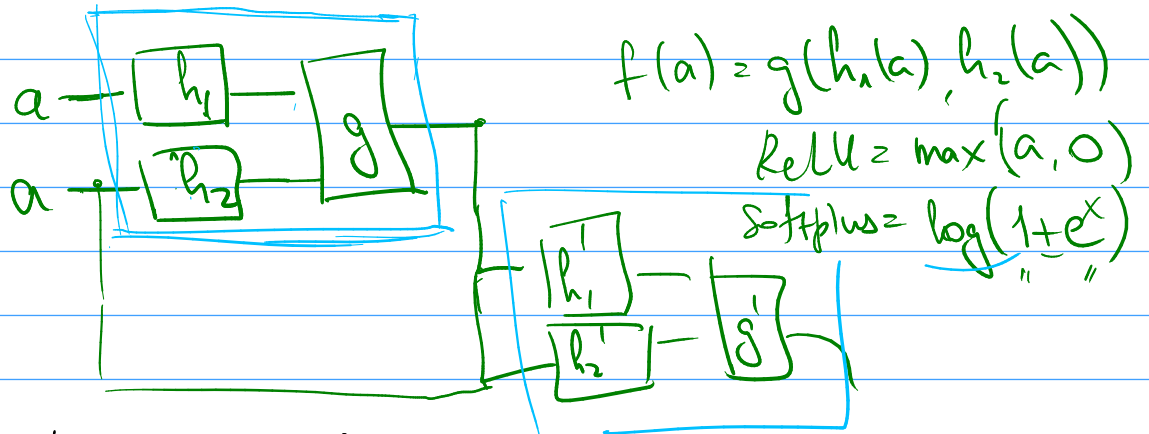


Neural architecture search  
Quoc Le

2017



$$f(a) = g(h_1(a), h_2(a))$$

$$\text{ReLU} = \max(a, 0)$$

$$\text{Softplus} = \log(1 + e^x)$$

$$\text{Swish}(a) = a \cdot \sigma(\beta \cdot a) = \frac{a}{1 + e^{-\beta a}}$$

Hard  Soft

$$\max(\dots)$$

$$\mathcal{S}_\beta(x_1, \dots, x_n) = \frac{\sum x_i \cdot e^{\beta x_i}}{\sum e^{\beta x_i}}$$

$$\mathcal{S}_\beta(g(x), h(x)) = g \cdot \frac{e^{\beta g}}{e^{\beta g} + e^{\beta h}} + h \cdot \frac{e^{\beta h}}{e^{\beta g} + e^{\beta h}} = \sigma(-x) = 1 - \sigma(x)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$= g \cdot \sigma(\beta(g-h)) + h \cdot \sigma(\beta(h-g))$$

$$= (g(x) - h(x)) \sigma(\beta \cdot (g(x) - h(x))) + h(x)$$

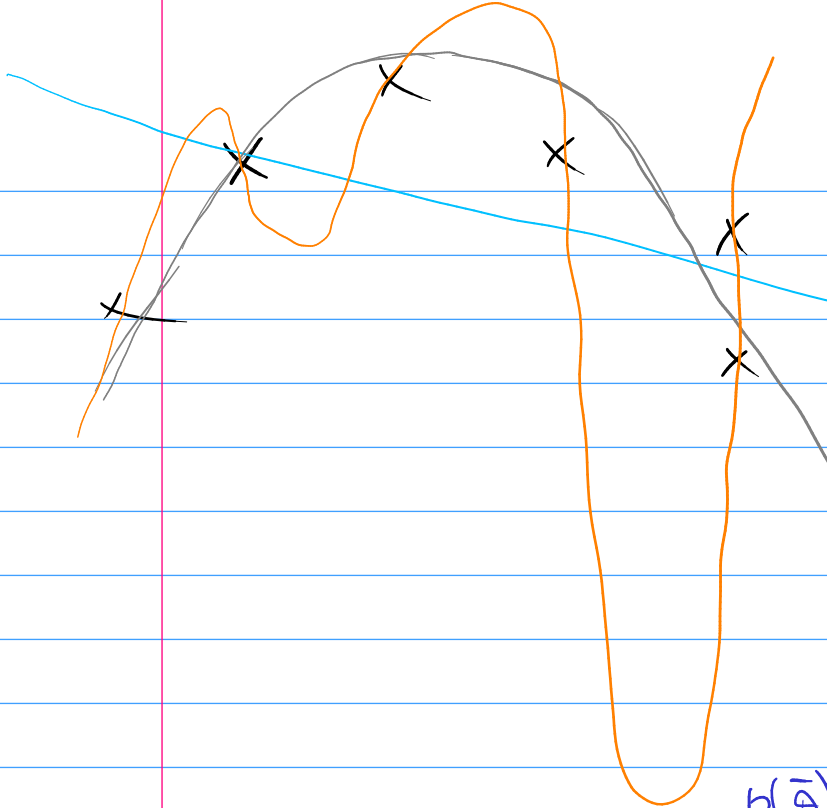
$$\mathcal{S}_\beta(x, 0) = x \cdot \sigma(\beta x)$$

$$\mathcal{S}_\beta(x, ax) = (1-a)x \cdot \sigma(\beta(1-a)x) + ax$$

$$\text{ACON} = \mathcal{S}_\beta(a_1 x, a_2 x)$$

$$\text{ReLU}(x) = \max(0, x)$$

$$\text{LReLU}(x) = \max(x, ax)$$



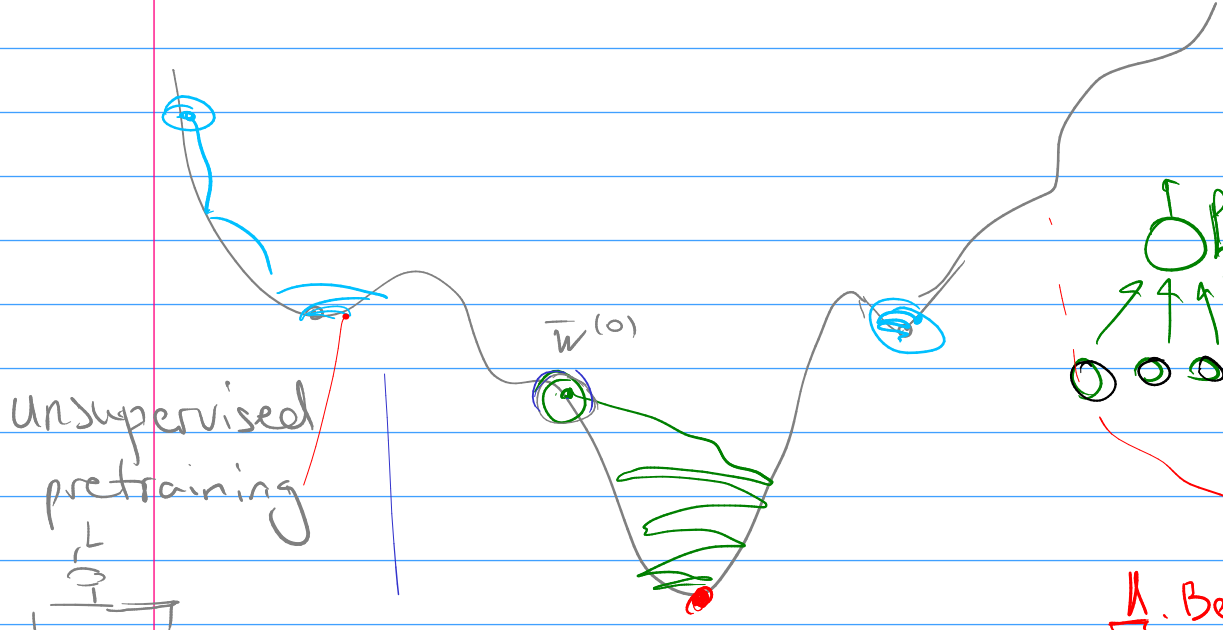
$$p(\theta|D) = \frac{p(\theta)p(D|\theta)}{p(D)}$$

posterior ←  $p(\theta|D)$   
 prior ←  $p(\theta)$   
 likelihood ←  $p(D|\theta)$

$$\log p(\theta|D) = \text{const} + \underbrace{\log p(D|\theta)} + \underbrace{\log p(\theta)}_{\text{regularizer}}$$

$$p(\bar{\theta}) = \mathcal{N}(\bar{\theta} | \bar{\theta}, \lambda^{-1}I) \quad \underbrace{-\frac{\lambda}{2} \|\bar{\theta}\|_2^2}_{L_2\text{-penalty}}$$

$(\bar{w})$   
 $L \rightarrow \min$

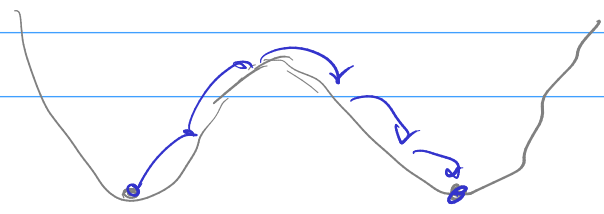
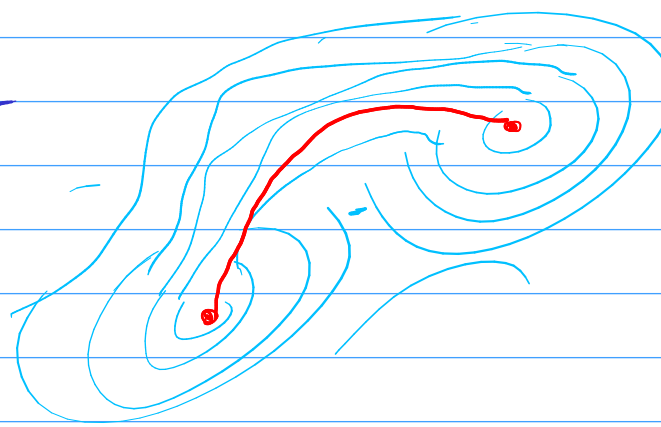
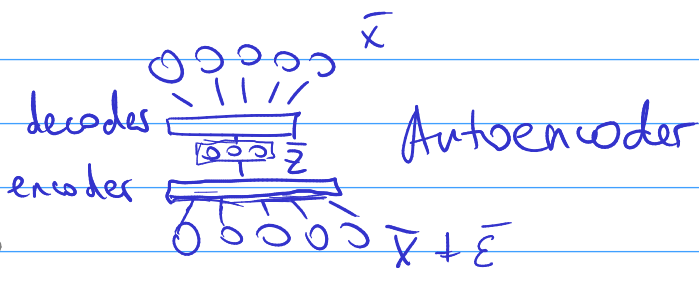
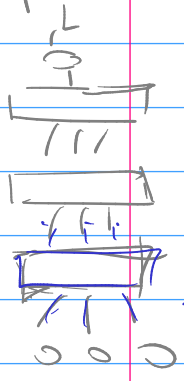


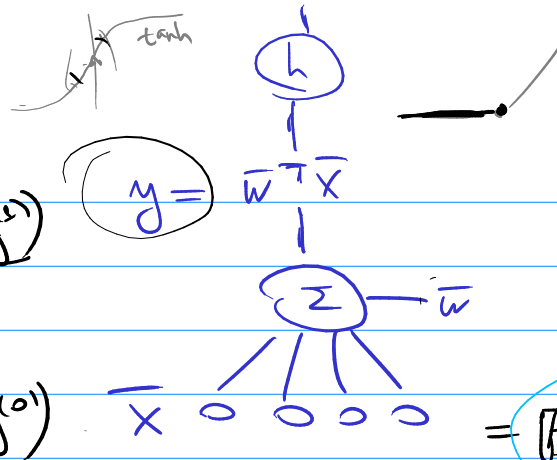
$$\text{Opt}(w = \bar{x})$$

$$(A^T A)^{-1} A^T y$$

1. Бессоб

unsupervised pretraining



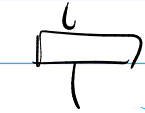


$$y = \overline{w^T x} = \sum \overline{w_i x_i} = \sum y_i$$

$$x^{(2)} = h(y^{(1)})$$



$$x^{(1)} = h(y^{(0)})$$



$$\overline{x} = \overline{x}^{(0)}$$

$$\begin{aligned} \text{Var}[y_i] &= \text{Var}[w_i x_i] = \\ &= E[w_i^2 x_i^2] - (E[w_i x_i])^2 \\ &= E[x_i^2] \text{Var}[w_i] + E[w_i]^2 \text{Var}[x_i] \\ &\quad + \text{Var}[x_i] \text{Var}[w_i] \end{aligned}$$

$$\text{Var}[y_i] = \text{Var}[x_i] \text{Var}[w_i]$$

$$\text{Var}[y] = \sum \text{Var}[y_i] = n \cdot \text{Var}[w_i] \cdot \text{Var}[x_i]$$

$$w_i \sim \mathcal{N}(w_i | 0, \frac{1}{3n})$$

$$w_i \sim \text{Unif}(\left[-\frac{\sqrt{3}}{\sqrt{n}}, \frac{\sqrt{3}}{\sqrt{n}}\right])$$

2010 Xavier Glorot

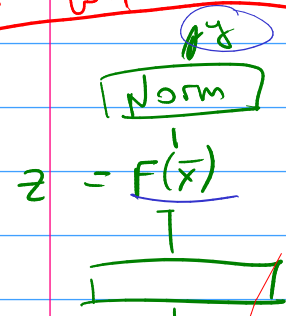
$$\text{Var}[a, b] = \frac{(b-a)^2}{12}$$

Xavier init

$$\text{Var}[y_i] = E[x_i^2] \text{Var}[w_i] + \text{Var}[x_i] \text{Var}[w_i]$$

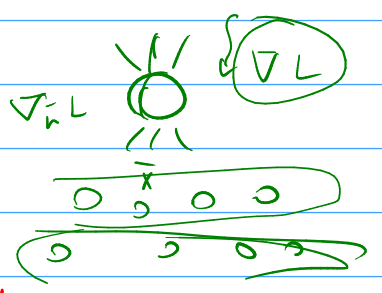
Kaiming He  
He init 2014

$$\text{Var } w_i = \frac{1}{3n}$$

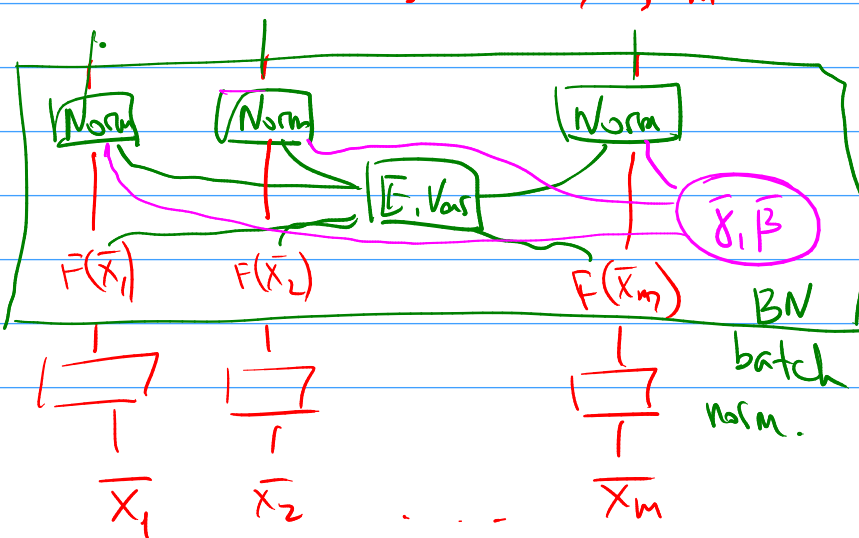
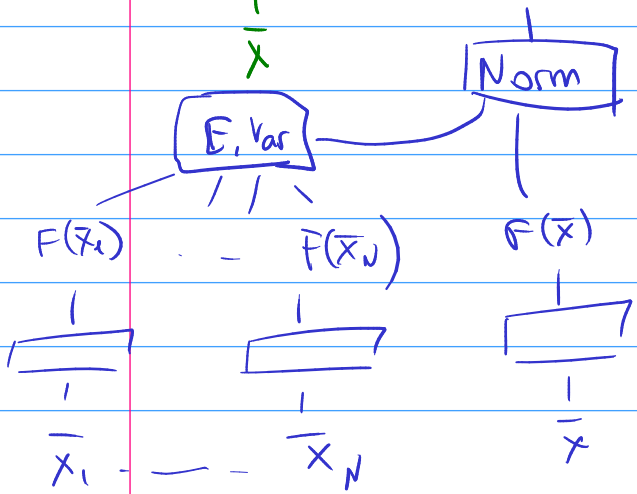


$$y_{jk} = \left( \frac{z_k - E[z_k]}{\text{Var}[z_k]} \right) \gamma_k + \beta_k$$

$$F(x) \rightarrow F(x) + \epsilon$$

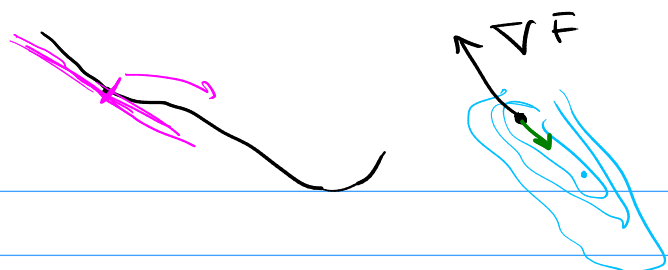


Mini-batches  $\overline{x}_1, \dots, \overline{x}_m$



BN  
batch norm.

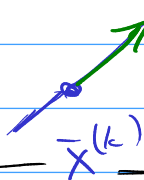
$$F(\bar{x}) \rightarrow \min$$



$$\bar{x}^{(k+1)} = \bar{x}^{(k)} - \alpha_k \cdot \nabla_{\bar{x}} F$$

$$\alpha_k = \arg \min F(\bar{x}^{(k)} - \alpha_k \nabla_{\bar{x}} F(\bar{x}^{(k)}))$$

Armijo, Wolfe



$$F(x, y) = x^2 + y^2$$

quasi-Newton

L-BFGS

holy  
grail

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} - \alpha_k \left( H_k^{-1} \nabla_{\bar{x}} F \right)$$

GD:  $\nabla_{\bar{w}} L = \nabla_{\bar{w}} \left[ \frac{1}{N} \sum_{n=1}^N \ell(\bar{x}_n, y_n, \bar{w}) \right] = \frac{1}{N} \sum_n \nabla_{\bar{w}} \ell(\bar{x}_n, y_n, \bar{w})$

SGD:

- $\bar{x}_1, \dots, \bar{x}_m \sim D$
- $\bar{w}^{(k+1)} = \bar{w}^{(k)} - \alpha_k \cdot \nabla L_m(\bar{w})$

$$F(\bar{w}) = \frac{1}{N} \sum_{i=1}^N f(\bar{x}_i, \bar{w}) = \mathbb{E} [f(\bar{x}, \bar{w})] \quad \bar{x} \sim \text{Unif}[D] \quad \frac{1}{m} \sum_{i=1}^m \nabla_{\bar{w}} \ell(\bar{x}_i, y_i, \bar{w})$$

$$\hat{F}_k(\bar{w}) = \frac{1}{m} \sum_{i=1}^m f(\bar{x}_i, \bar{w}) \approx \mathbb{E} [f(\bar{x}, \bar{w})]$$

$$\nabla_{\bar{w}} \hat{F}_k(\bar{w}) \approx \nabla_{\bar{w}} F(\bar{w})$$