# Deep Networks

Overfitting + Generalization

p parameters, n datapoints, d dimensions

$p \approx d$. ← Classic

$p \gg nd$ ← Modern NNs

$p \approx 1$ trillion    $n = 60000, d = 28 \times 28$

Avoiding Overfitting

- Bagging
- Regularization

$\hookrightarrow \nabla_\theta \quad L(x, y, \theta) + \frac{\lambda}{2} \|\theta\|_2^2$

$\Rightarrow \nabla_\theta L(x, y, \theta_{t-1}) + \lambda \theta_{t-1}$

$\theta_t = \theta_{t-1} - \eta \left( \nabla_\theta L(x, y, \theta_{t-1}) + \lambda \theta_{t-1} \right)$

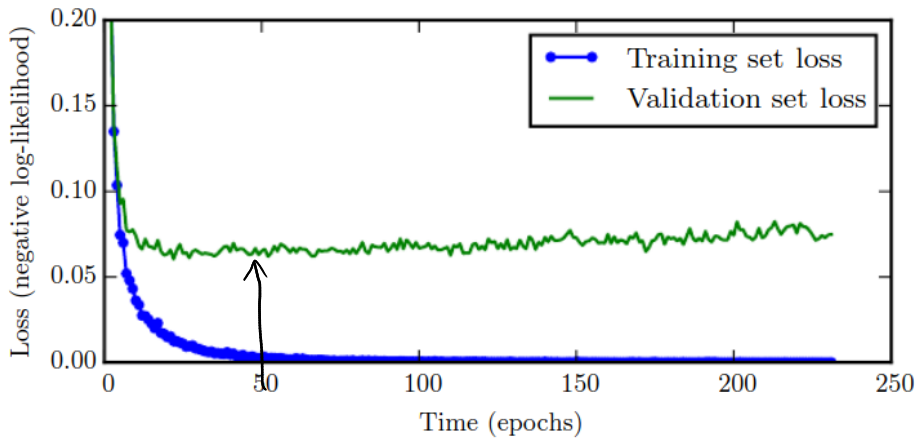$\theta_t = (1 - \eta\lambda)\theta_{t-1} - \eta \nabla_\theta L(x, y, \theta_{t-1})$

$\quad \quad \hookleftarrow$ Weight decay

- Data Augmentation

$$6 \rightarrow \bcancel{\phi} \rightarrow \bcancel{\phi}$$

# Early Stopping



Checkpointing

## Dropout
Keep each node w.p. $p > 0$, when training each point

$p = 0.5$



$x_1$
$x_2$
$x_3$
$x_4$

$\hookrightarrow$ $x_1$ $x_2$ $x_3$ $x_4$



$h^{(1)}$

$\bigcirc$ — $\frac{1}{0.5} \cdot 1 \cdot h_1^{(1)}$

$\bigotimes$ — $\frac{1}{0.5} \cdot 0 \cdot h_2^{(1)}$

$\bigcirc$ — $\frac{1}{0.5} \cdot 1 \cdot h_3^{(1)}$

$\bigotimes$ — $\frac{1}{0.5} \cdot 0 \cdot h_4^{(1)}$

$$z^{(2)} = \underbrace{W^{(2)} h^{(1)}}_{\approx p \text{ times its prev value}} + b^{(2)}$$

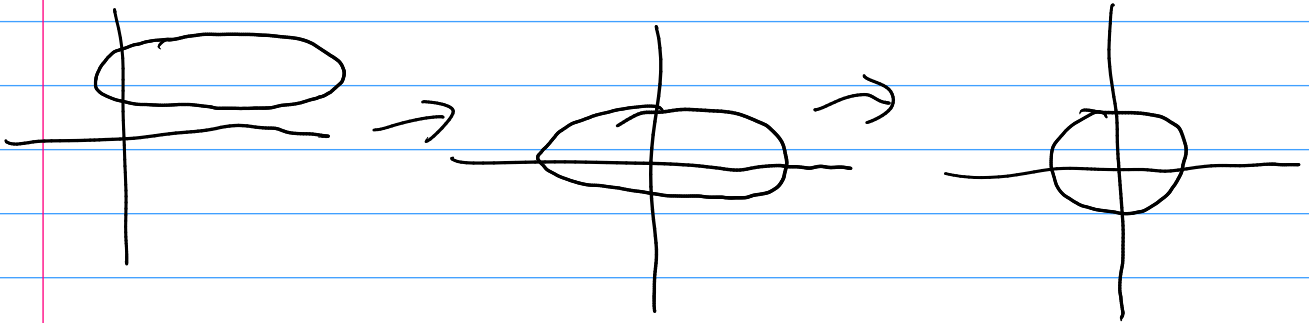Multiply $h^{(1)}$ by $\frac{1}{p}$. (before Dropout)

(Inverted dropout)

Test time: No dropout, no scaling

# Normalization

Normalize features before training

$\mu = \frac{1}{n}\sum X_i$. $X_i = X_i - \mu$. $\sigma_j^2 = \frac{1}{n}\sum X_{ij}^2$. $X_{ij} = X_{ij}/\sigma_j$

$\uparrow \forall i$ $\qquad\qquad\qquad\qquad\qquad \uparrow \forall i,j$



# Batch norm

$z/h$



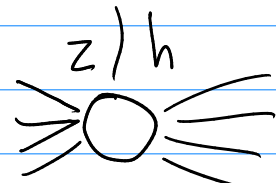| | |
|---|---|
| **Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$; Parameters to be learned: $\gamma, \beta$ | |
| **Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$ | |
| $\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i$ | // mini-batch mean |
| $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2$ | // mini-batch variance |
| $\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ | // normalize |
| $y_i \leftarrow \gamma\hat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$ | // scale and shift |

learnable params
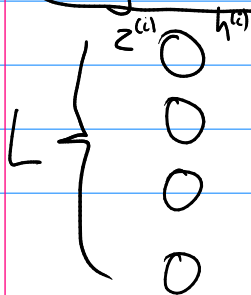
$z^{(i)} = W h^{(i-1)} + b^{(i)}$

$h^{(i)} = f(z^{(i)})$

Minibatch SGD
- Process 64 pts at once.
(128, 256)

Applied to each neuron individually
Averaging over a minibatch.

## Layer norm

$z^{(i)} \quad h^{(i)}$

$L \{$ (neurons)

For each datapoint:

$\mu^{(i)} = \frac{1}{L}\sum_{j=1}^{L} z_j^{(i)} \qquad z_j^{(i)} = z_j^{(i)} - \mu^{(i)}$

$\sigma^{(i)^2} = \frac{1}{L}\sum z_j^{(i)^2} \qquad z_j^{(i)} = \frac{z_j^{(i)}}{\sigma^{(i)^2}}$

Averaged over a layer

# Optimization

First-order methods: 1st derivative
Second-order " " 2nd " "

## Gradient Descent    "Batch"

$$\theta \leftarrow \theta - \eta \cdot \underbrace{\frac{1}{n} \sum_{i=1}^{n} \nabla_\theta l(x_i, y_i, \theta)}_{} \qquad \mathbb{E}_{(x,y) \sim D}\left[\nabla_\theta l(x,y,\theta)\right]$$

↖ hyper parameter

## Stochastic GD

$$\theta \leftarrow \theta - \eta \cdot \nabla_\theta l(x_i, y_i, \theta)$$

Shuffling after each epoch

## Mini-batch SGD

$$\theta \leftarrow \theta - \eta \cdot \frac{1}{|B|} \cdot \sum_{i \in B} \nabla_\theta l(x_i, y_i, \theta)$$

$$\underbrace{x_1, x_2,}_{} \underbrace{x_3, \ldots,}_{} \underbrace{x_n}_{}$$

$|B| \approx 64$ to $256$

## Challenges

- $\eta$?
- LR schedules ← doesn't "adapt" to data
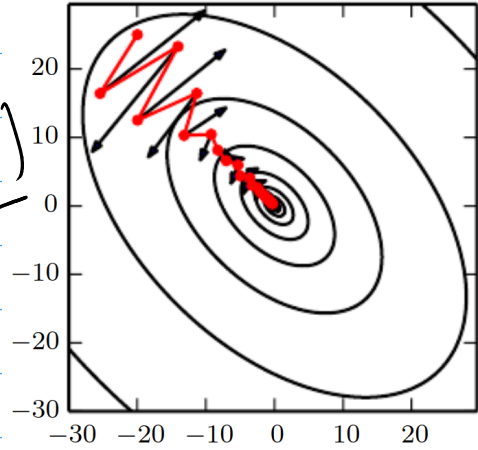- Can we have diff $\eta$ for diff coords?

"Momentum"  $\gamma < 1$.
$$\gamma = 0.9?$$
$$V_t = \gamma V_{t-1} + (1-\gamma)\eta \underbrace{\frac{1}{|B|}\sum \nabla \ell_{(x_i, y_i, \theta)}}_{g_t}$$

$$\theta \leftarrow \theta - V_t$$

$$V_t = 0.1 g_t + 0.1 \cdot 0.9 g_{t-1} + 0.1 \cdot 0.9^2 g_{t-2}$$
$$+ \cdots + 0.1 \cdot 0.9^{t-1} g_1$$

Total coeff: $1-\gamma^t \approx 1$ when $t$ large



Black: Gradient
Red: momentum + SGD

# Nesterov Momentum/Accelerated Gradient

$$V_t = \gamma V_{t-1} + (1-\gamma)\eta \frac{1}{|B|}\sum \nabla \ell_{(x_i, y_i, \theta - \gamma V_{t-1})}$$
$$\theta = \theta - V_t$$

## Adaptive Learning Rates
Change LR over algo.
Based on "importance" of each param.
- Lot of updates → low LR for param
- Few updates → large LR for param

$g_t \in \mathbb{R}^P$ is gradient at time $t$.

SGD → $\theta_{t+1, i} = \theta_{t,i} - \eta g_{t,i}$ ← coord $i$ update

Define $G_{t,i} = \sum_{j=1}^{t} g_{j,i}^2$ (sum of squared $\overset{grad}{updates}$)

AdaGrad

$\rightarrow \theta_{t+1,i} = \theta_{t,i} - \dfrac{\eta}{\sqrt{G_{t,i} + \varepsilon}} g_{t,i}$

$\underbrace{\qquad}_{} \leftarrow \approx 10^{-8}$

RMSProp: $G_{t,i} = 0.9 G_{t-1,i} + 0.1 g_{t,i}^2$

Momentum, but on 2nd moment (grad$^2$) rather than 1st moment

$\theta_{t+1,i} = \theta_{t,i} - \dfrac{\eta}{\sqrt{G_{t,i} + \varepsilon}} g_{t,i}$

## Adam

- Momentum + RMSProp (plus bias correction)
- $\beta_1, \beta_2$ hyper params

$m_{t,i} = \beta_1 m_{t-1,i} + (1-\beta_1) g_{t,i} \quad \longleftarrow$ Momentum

$V_{t,i} = \beta_2 V_{t-1,i} + (1-\beta_2) g_{t,i}^2 \quad \leftarrow$ RMS Prop updates

$\hat{m}_{t,i} = \dfrac{m_{t,i}}{1-\beta_1^t} , \quad \hat{V}_{t,i} = \dfrac{V_{t,i}}{1-\beta_2^t} \quad \leftarrow$ Bias correction

$\theta_{t+1,i} = \theta_t - \dfrac{\eta}{\sqrt{\hat{V}_{t,i}} + \varepsilon} \hat{m}_{t,i}$

$\beta_1 = 0.9, \ \beta_2 = 0.999, \ \varepsilon = 10^{-8}$