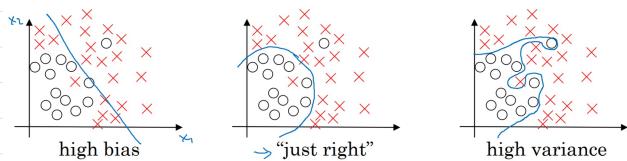


Train / Dev / Test sets

data	train	(dev) test
Normal	60% 6'000	20% 20% 2'000 2'000
Big Data	99.5% 10'000'000	0.25% 0.25% 10'000 10'000

→ Hold-out Cross Validation
 → Development set

* dev & test set come from same distribution.



Train error:	1%	15%	16%
Dev error:	11%	16%	30%
	high v	high b	high v&b

Fixes →

High bias:
 - Bigger Network
 - Train Longer

High variance:
 - More Data
 - Regularization

* Deep learning does not have the "b-v trade off"

Regularization

$$J(w, b) = \frac{1}{m} \sum_i^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + L_1 + L_2$$

logistic Reg →

$$L_1 = \frac{\lambda}{2m} \|w\|_1 = \frac{\lambda}{2m} \sum_j^n |w_j|$$

$$L_2 = \frac{\lambda}{2m} \|w\|_2^2 = \frac{\lambda}{2m} \sum_j^n w_j^2 = \frac{\lambda}{2m} w^T w$$

NN → $L_2 = \frac{\lambda}{2m} \sum_i^L \|w^{[l]}\|_F = \frac{\lambda}{2m} \sum_i^L \sum_j^n \sum_{ij}^{[l]} (w_{ij}^{[l]})^2$

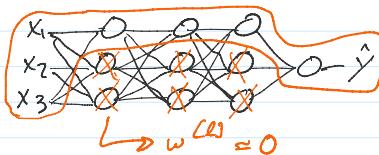
↳ Frobenius norm

$$\Rightarrow dw^{[l]} = (\text{back prop}) + \frac{\lambda}{m} w^{[l]}$$

$$\rightarrow w^{[l]} = \alpha dw^{[l]}$$

weigh decay → $w^{[l]} = (1 - \frac{\alpha \lambda}{m}) w^{[l]} - \alpha (\text{back prop})$

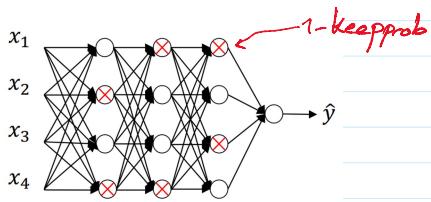
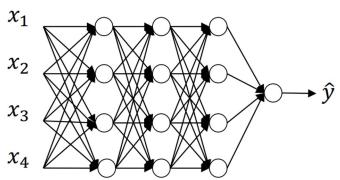
$$\uparrow \downarrow \quad w^{[l]} \downarrow z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$



→ makes layers linear



Dropout Regularization



$$d^3 = \text{np.random}(size=a^3) < \text{keepprob}$$

$$a^3 = \text{np.multiply}(a^3, d^3)$$

$$a^3 /= \text{keepprob}$$

Inverted Dropout

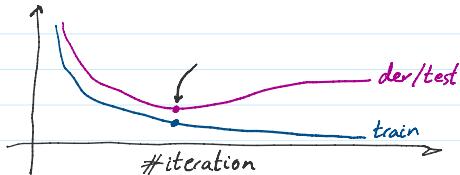
*not during test

*Can't rely on any One Feature

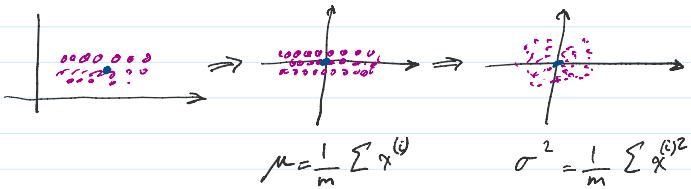
Data Augmentation

$$4 \rightarrow \begin{matrix} 4 \\ 4 \\ 4 \\ 4 \end{matrix}$$

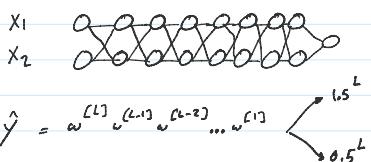
Early Stopping



Normalizing Training Sets



Vanishing / Exploding gradients problem



Fix:

$$w^{[l]} = \text{np.random()} * \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

$$* \quad \left(\frac{1}{n^{[l-1]}}\right) \quad \text{Xavier}$$

$$* \quad \left(\frac{2}{n^{[l-1]} + n^{[l]}}\right)$$

Mini-Batch

1

$\frac{1}{M_M}$

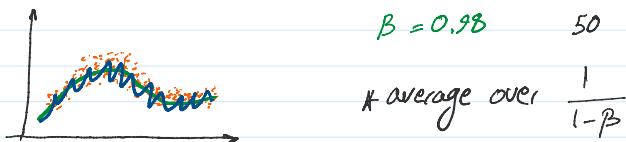
st



- * size = m \rightarrow Batch GD Too long per filter
- * size = 1 \rightarrow Stochastic GD loose speed from vec-
- * size is usually $2^k = 64, 128, 256, \dots$

Exponentially Weighted Averages

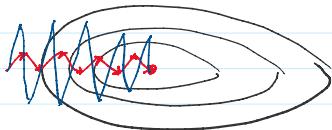
$$v_t = \beta v_{t-1} + (1-\beta) \theta_t \quad \beta = 0.9$$



Bias Correction

$$\rightarrow v_t = \frac{v_t}{1-\beta^t}$$

Gradient Descent with Momentum



On #iter t:

$$v_{dw} = \beta v_{dw} + (1-\beta) dw$$

$$v_{db} = \beta v_{db} + (1-\beta) db$$

$$w_- = \alpha v_{dw} \quad b_- = \alpha v_{db}$$

RMS prop

On #iter t:

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2$$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

$$w_- = \frac{\alpha dw}{\sqrt{S_{dw} + \epsilon}} \quad b_- = \frac{\alpha db}{\sqrt{S_{db} + \epsilon}}$$

Adam Optimization Algorithm

On #iter t:

$$v_{dw} = \beta_1 v_{dw} + (1-\beta_1) dw \quad v_{db} =$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \quad S_{db} =$$

$$v_{dw}^{corr} = \frac{v_{dw}}{1-\beta_1^t}$$

α needs to be tune

$$V_{dw}^{corr} = \frac{V_{dw}}{1 - \beta_1 t}$$

α needs to be tune

$$S_{dw}^{corr} = \frac{S_{dw}}{1 - \beta_2 t}$$

$$\beta_1 = 0.9$$

$$\omega_- = \frac{\alpha V_{dw}^{corr}}{\sqrt{S_{dw}^{corr} + \epsilon}}$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

$$b_- = \frac{\alpha V_{db}^{corr}}{\sqrt{S_{db}^{corr} + \epsilon}}$$

Learning Rate Decay

$$\textcircled{1} \quad \alpha = \frac{1}{1 + \text{decay} * \text{epoch}} \alpha_0$$

$$\textcircled{2} \quad \alpha = 0.95^{\text{epoch}} \alpha_0$$

$$\textcircled{3} \quad \alpha = \frac{K}{\text{epoch}} \alpha_0$$

α  staircase

Tuning Process

1 - α

2 - β

- hidden layers

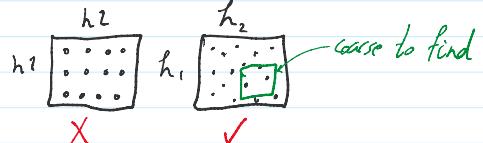
- mini batch size

3 - layers

- learning rate decay

4 - $\beta_1, \beta_2, \epsilon$

* Use random values for hyperparameters



* babysitting (Panda) vs many in parallel (Carrot)