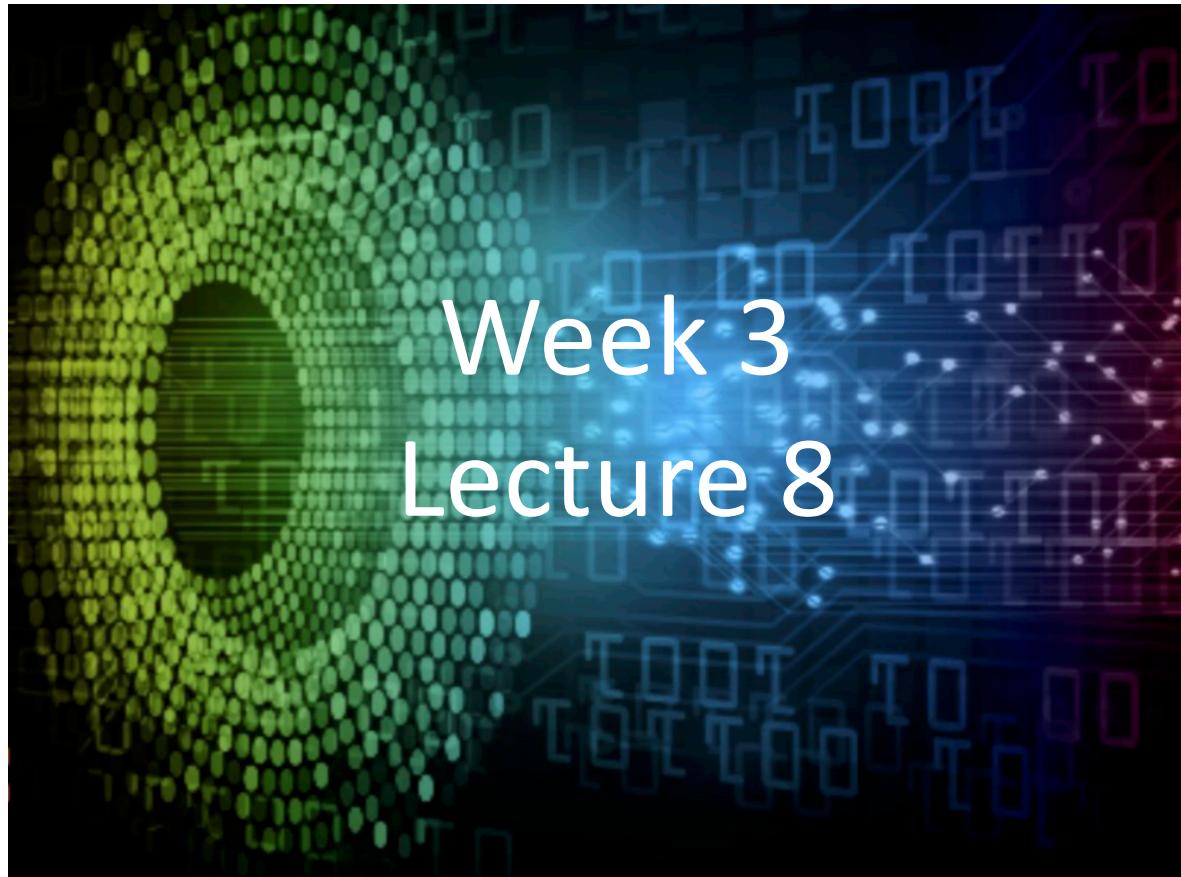


Introduction to Deep Learning Applications and Theory



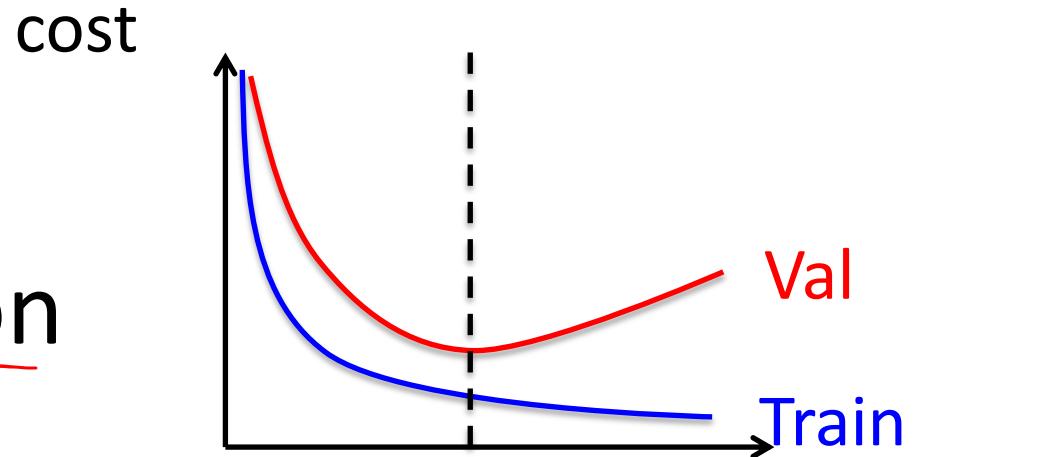
ECE 596 / AMATH 563

Previous Lecture: Optimization of Stochastic Data

1. Extensions of Gradient Descent for Stochastic Data

- *optimization*

2. Cross Validation

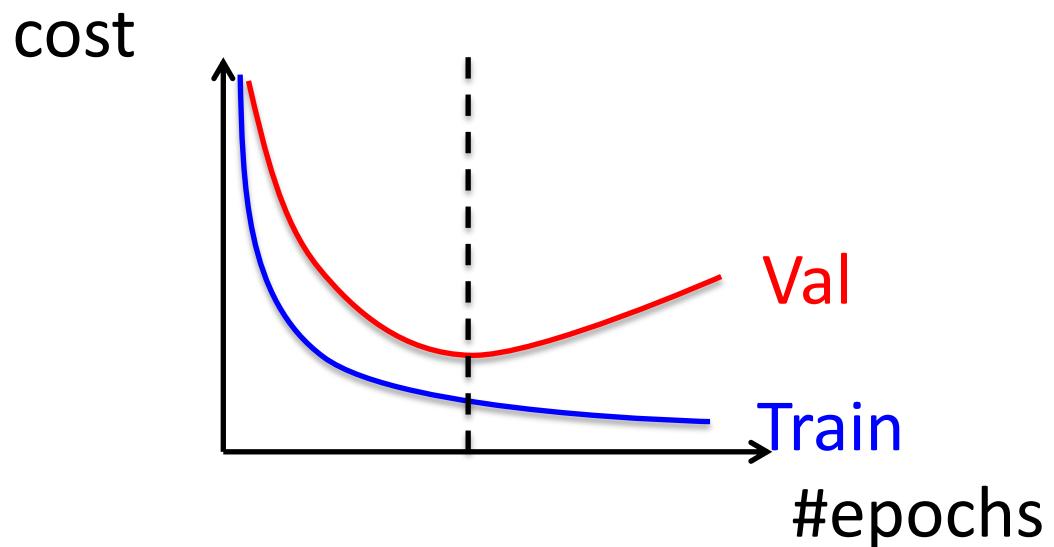


- *Data*
=> *Overfitting*
Under

Current Lecture: Regularization

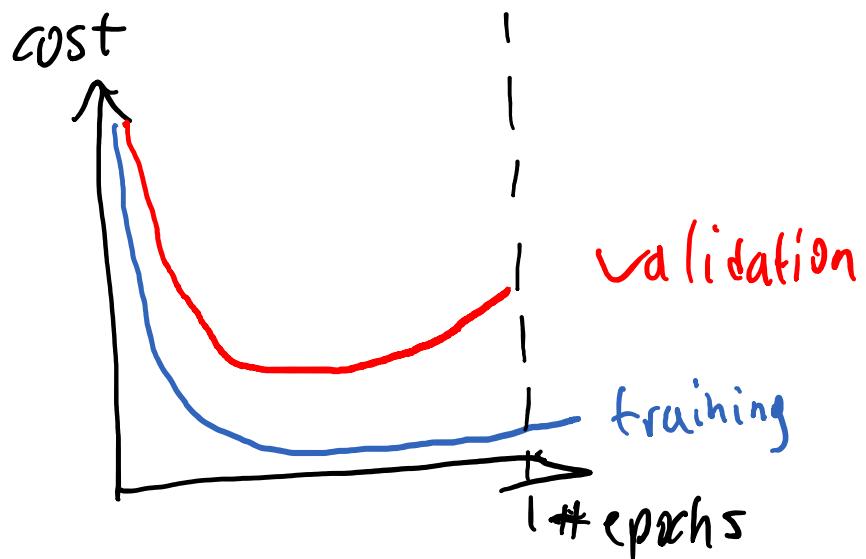
1. Regularization

- L2/L1
- Dropout
- Additional



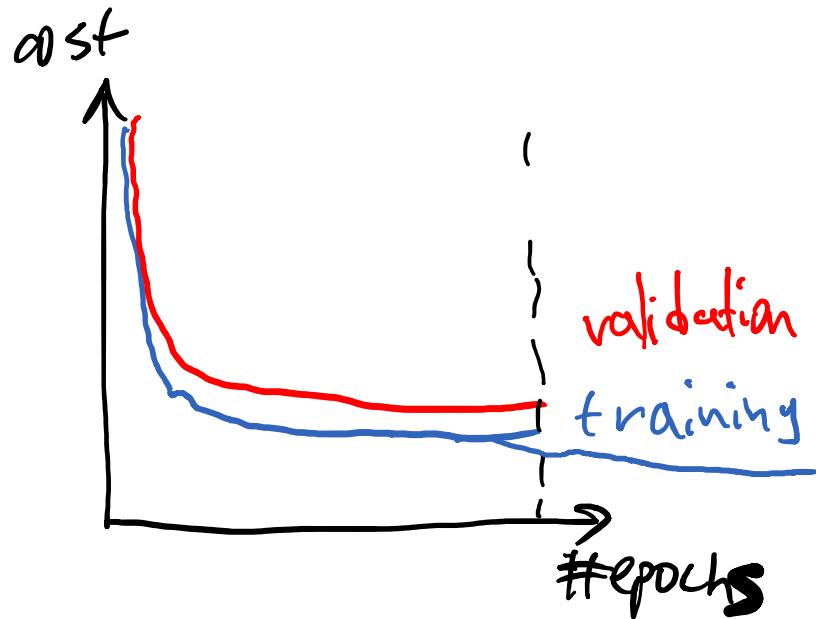
High Variance vs High Bias

Overfitting



High-variance

Underfitting



High-bias

High Variance vs High Bias

Overfitting

"Not enough data "for parameters"

More Data

Regularization

Dropout

Initialization

Optimization
solutions

Underfitting

"Not enough parameters"

More Layers/Nodes

Longer Training

Architecture

Deep vs CNN

Hyperparam optim

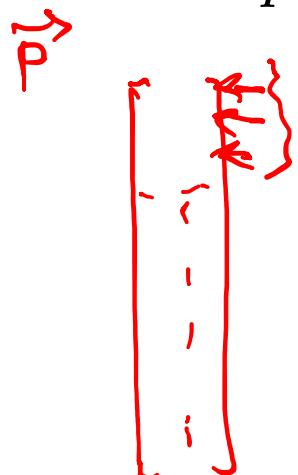
Over- and Under-Determined

$$\vec{p}^* = \arg \min_{\vec{p}} L(\vec{p})$$
$$\underline{A\vec{p} = \vec{b}}$$

Over-determined:

$$\arg \min_{\vec{p}} \left(\overline{\|A\vec{p} - \vec{b}\|_2} + \lambda g(p) \right)$$

Objective *regularizer*

$$\vec{p} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$


$$g(p) = \|p\|_2$$
$$= \|p\|_1$$

Regularization

Add regularization term to the cost

L2 or L1

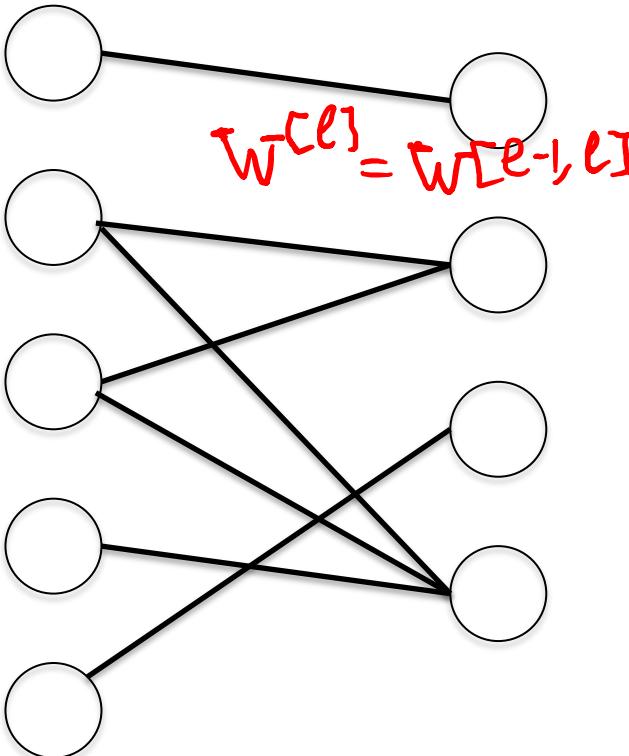
For single layer

$$J(\vec{w}, b) = \underbrace{\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})}_{J} + \underbrace{\frac{\lambda}{2m} \|\vec{w}\|_2^2}_{R_2} + \underbrace{\frac{\lambda}{m} \|\vec{w}\|_1}_{R_1}$$

MultiLayered Regularization

Layer: $\ell-1$

$$\vec{z}^{[\ell-1]} \rightarrow \vec{a}^{[\ell-1]}$$



$$\#n^{\ell-1}$$

$$\#n^\ell$$



$$\vec{z}^{[\ell]} \rightarrow \vec{a}^{[\ell]}$$

$$\vec{z}^{[\ell]} = W^{[\ell]} \cdot \vec{a}^{[\ell-1]} + b^{[\ell]}$$

$$\vec{a}^{[\ell]} = f(\vec{z}^{[\ell]})$$

Multilayered Regularization

$$J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = \underbrace{\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})}_{\text{J}} + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^L ||W^{[l]}||_F^2}_{\text{R}}$$

$$\underbrace{||W^{[l]}||_F^2}_{\text{Frobenius}} = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2$$

Regularization effect in GD

$$w_{k+1}^{[l]} = \underline{w_k^{[l]}} - \alpha \nabla_{w_k^{[l]}} J^R$$

GD

$$\nabla_{w_k^{[l]}} J^R = \nabla_{w_k^{[l]}} J + \frac{\lambda}{m} w_k^{[l]}$$

not additional

$$w_{k+1}^{[l]} = \left(1 - \frac{\alpha \lambda}{m}\right) w_k^{[l]} - \alpha \nabla_{w_k^{[l]}} J$$

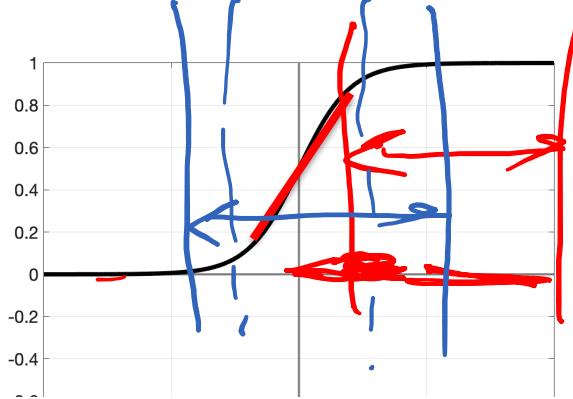
damping

L2 Regularization Intuition

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|\vec{w}\|_2^2$$

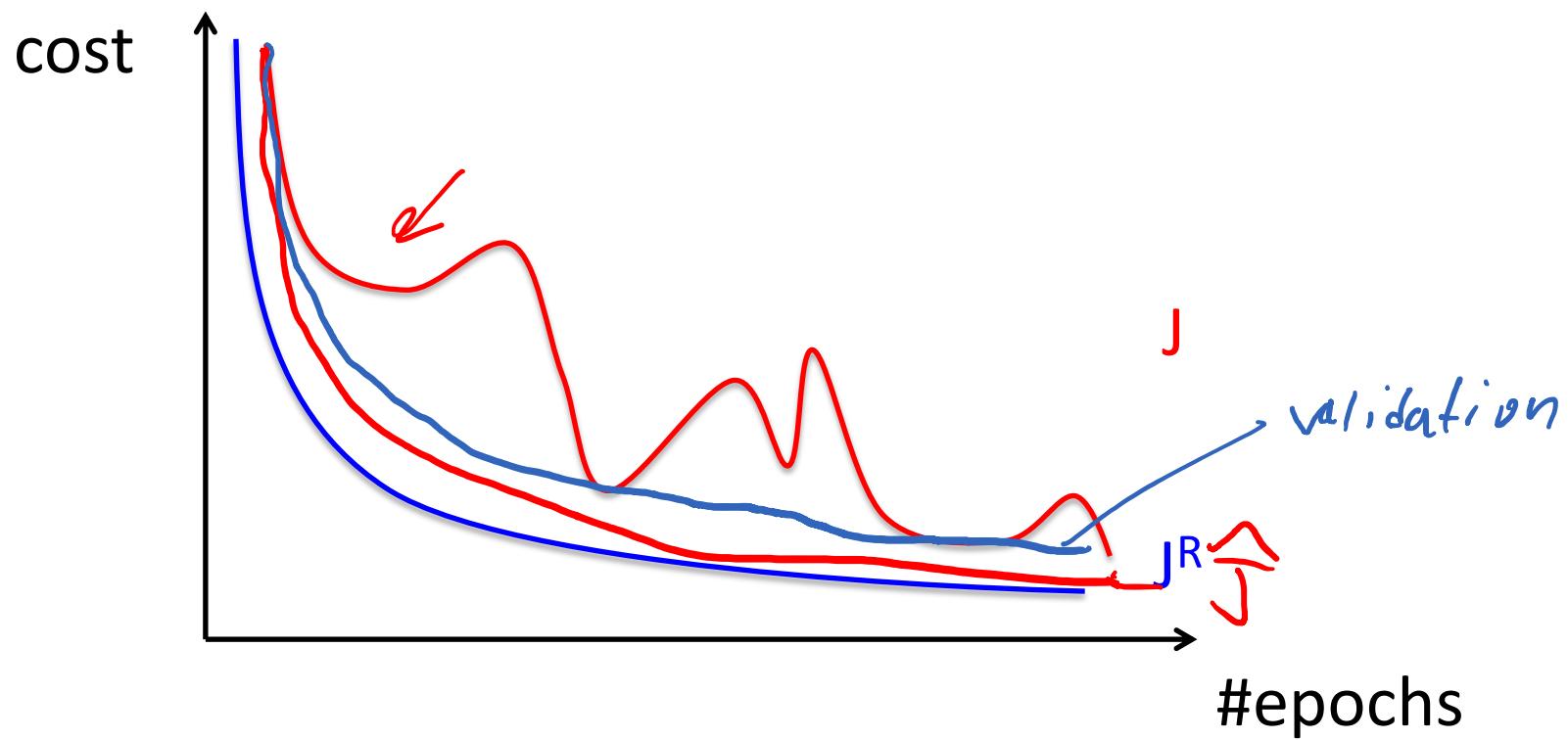
$\lambda \rightarrow \text{large}$

$w^{[l]} \approx 0$

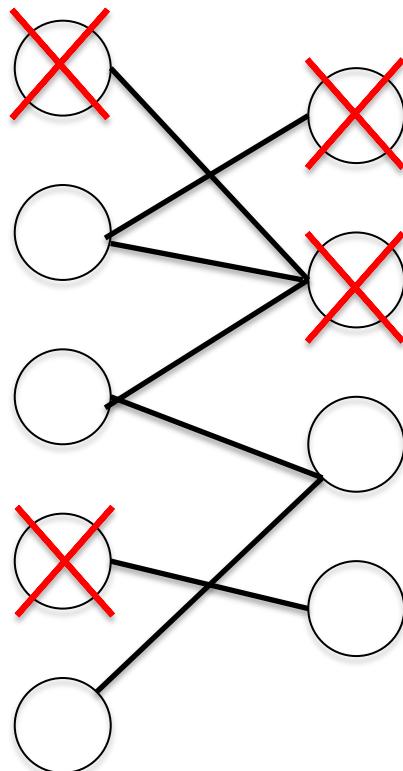


Have to be careful
that $\vec{w} \neq 0$.

Cost curve



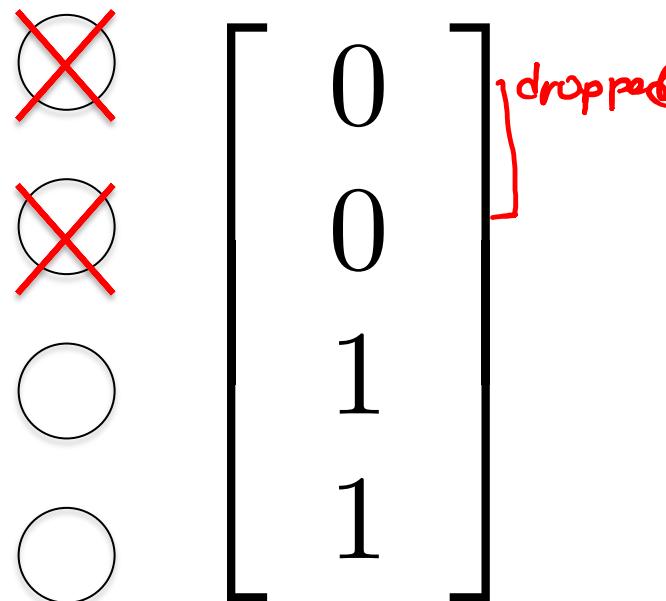
Dropout Regularization



$$\underline{P_d = 0.5}$$

Dropout - disabling activation

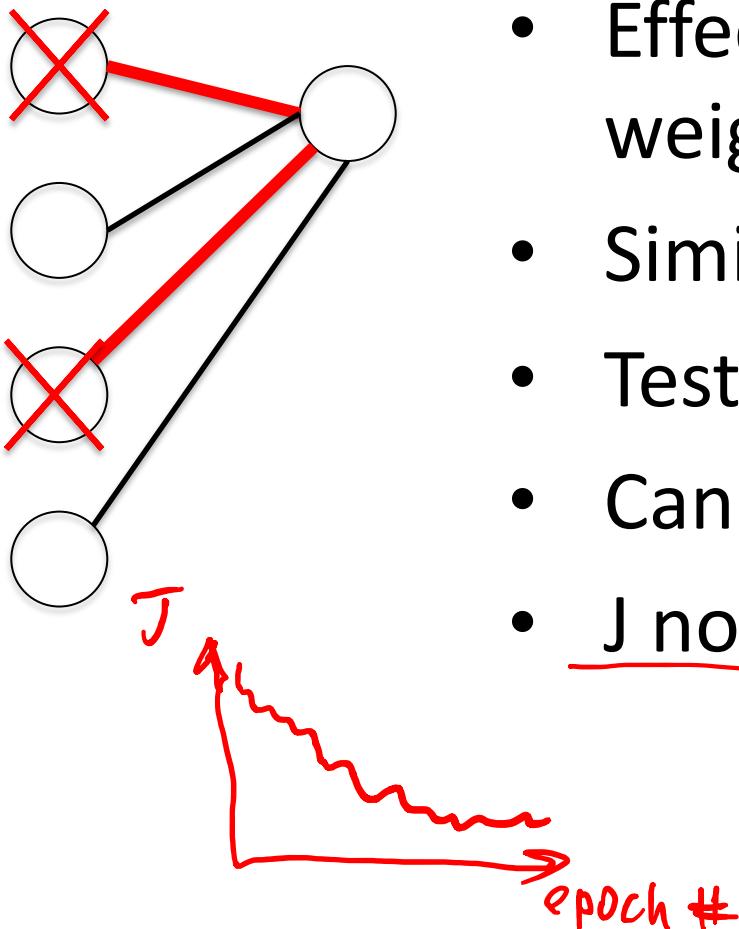
Inverted Dropout



$$\underline{\vec{a}_l^d} = \underline{\vec{a}_l} \times \underline{\vec{d}_l} / (\underline{1 - p_d})$$

Dropout intuition

"restricting optimization to fewer params
at a time"

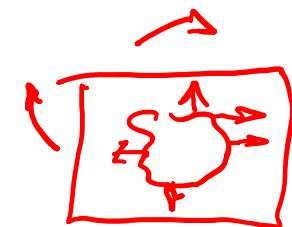


- Effectively spreading the weights
- Similar to L2 reg
- Testing with dropout $p_d=0$
- Can be W dim dependent
- J not well defined

Additional Regularization

- Data Augmentation

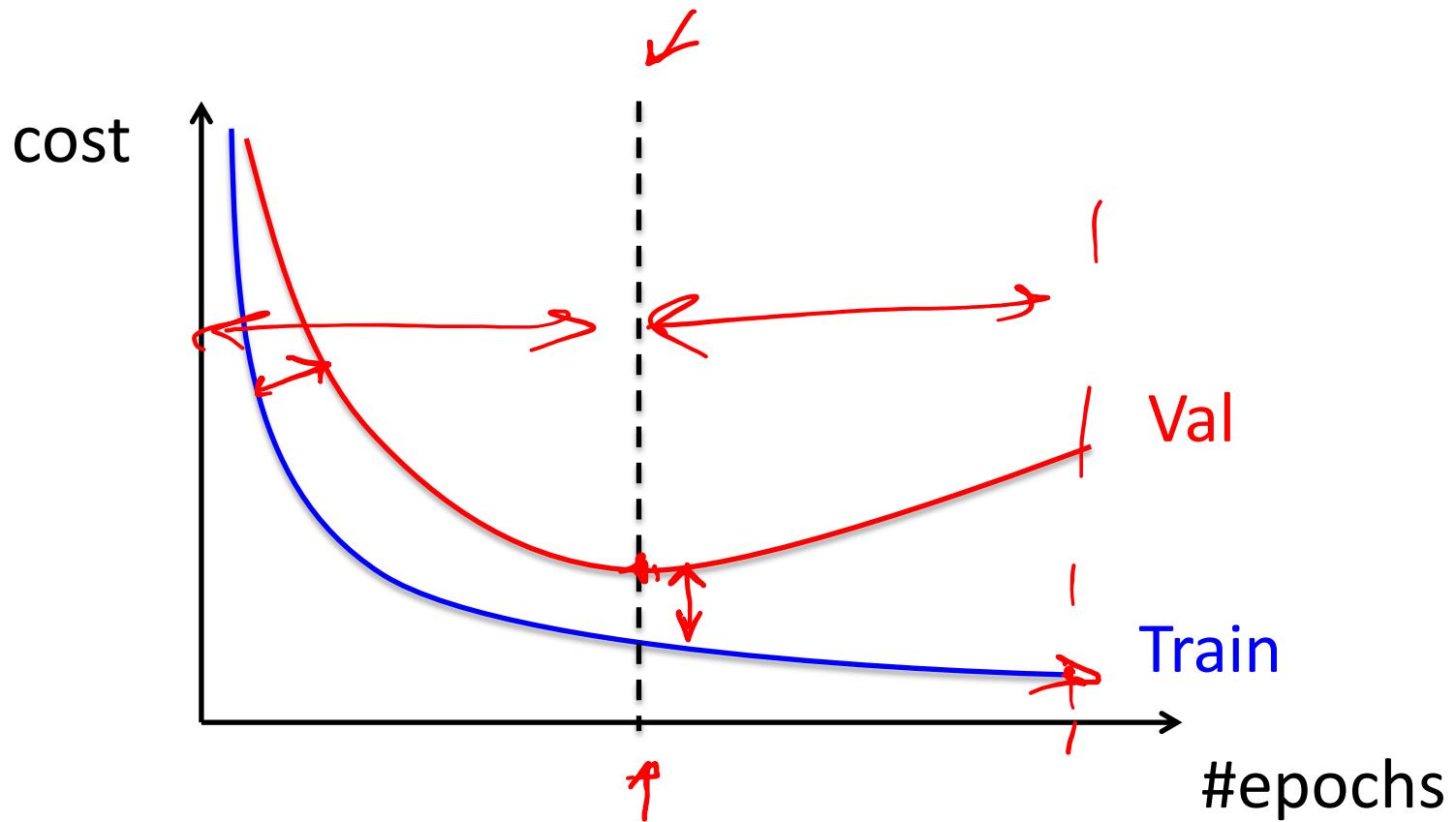
Enhance the data set with additional manipulations



General input: add noise/distortions, synthetic
Images: resolutions, rotate, add symmetries
Shapes/digits: distort

Additional Regularization

- Early Stopping:



Balance

Optimization of J

GD



Over-fitting

Regularization