# Regularization

Seyoung Yun

- http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

- N. Srivstava et al, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting" http://jmlr.org/papers/volume15/srivastava14a/srivastava14a.pdf

- Sergey Ioffe and Christian Szegedy"Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift" https://arxiv.org/abs/1502.03167

- C. Zhang et al "Understanding deep learning requires rethinking generalization" https://arxiv.org/abs/1611.03530

# Regularization

*training data*

*model*

*loss function*

*optimization*

- **"A regularizer is anything that hurts the training process"**
  C. Zhung at ICLR2017 ([https://www.youtube.com/watch?v=kCj51pTQPKI](https://www.youtube.com/watch?v=kCj51pTQPKI))

  - data augmentation

  - **weight decay - with an additional cost**

  - **dropout - by adding random noise**

*bias & variance.*

# Linear Regression

- RSS: cost of linear regression

$$\mathcal{L}(w,b) = \sum_{i=1}^{m} (y^{(i)} - w^\top x^{(i)} - b)^2$$

_linear function_

- regularizer

$$\mathcal{L}(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y^{(i)} - w^\top x^{(i)} - b)^2 + \frac{\lambda}{2m} \|w\|_2^2$$

or

$$\mathcal{L}(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y^{(i)} - w^\top x^{(i)} - b)^2 + \frac{\lambda}{2m} \|w\|_1$$

# Weight Decay

$$J(W^{[1]}, b^{[1]}, \cdots, W^{[L]}, b^{[L]})$$

$$= \frac{1}{m} \sum_{i=1}^{m} \ell\left(\hat{Y}^{(i)}, Y^{(i)}\right) + \frac{\lambda}{2m} \sum_{l=1}^{L} \| W^{[l]} \|_F^2$$
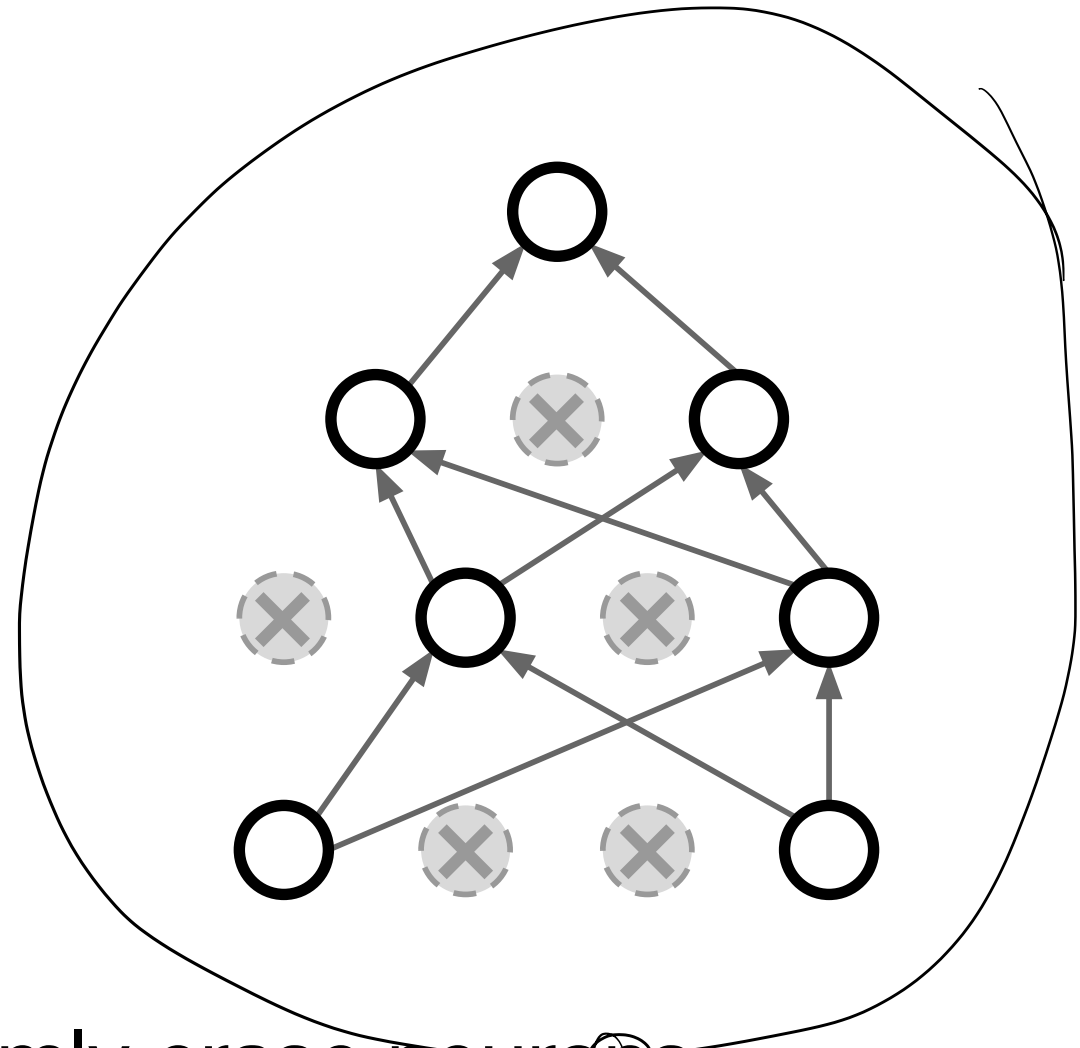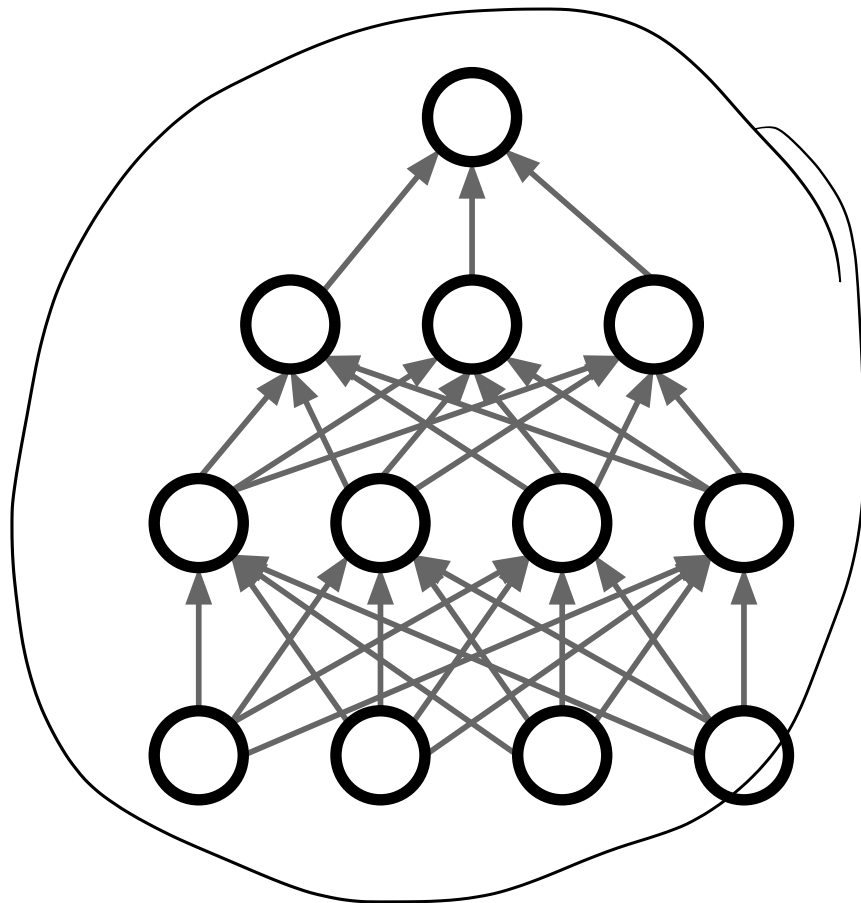
$$\| W^{[l]} \|_F^2 = \sum \sum \left( W_{jk}^{[l]} \right)^2$$

$$\Rightarrow \nabla_{W^{[l]}} J = \underbrace{\frac{1}{m} \sum_{i=1}^{m} \nabla_{W^{[l]}} \ell\left(\hat{Y}^{(i)}, Y^{(i)}\right)}_{\text{Back prop.}} + \frac{\lambda}{m} \cdot W^{[l]}$$

$$W^{[l]}(t+1) = W^{[l]}(t) - \gamma \cdot \nabla_{W^{[l]}} J\left(W^{[l]}(t)\right)$$

$$= \left(1 - \frac{\gamma\lambda}{m}\right) W^{[l]}(t) - \underbrace{\frac{\gamma}{m} \sum_{i=1}^{m} \nabla_{W^{[l]}} \ell\left(\hat{Y}^{(i)}, Y^{(i)}\right)}_{\text{Back prop.}}$$
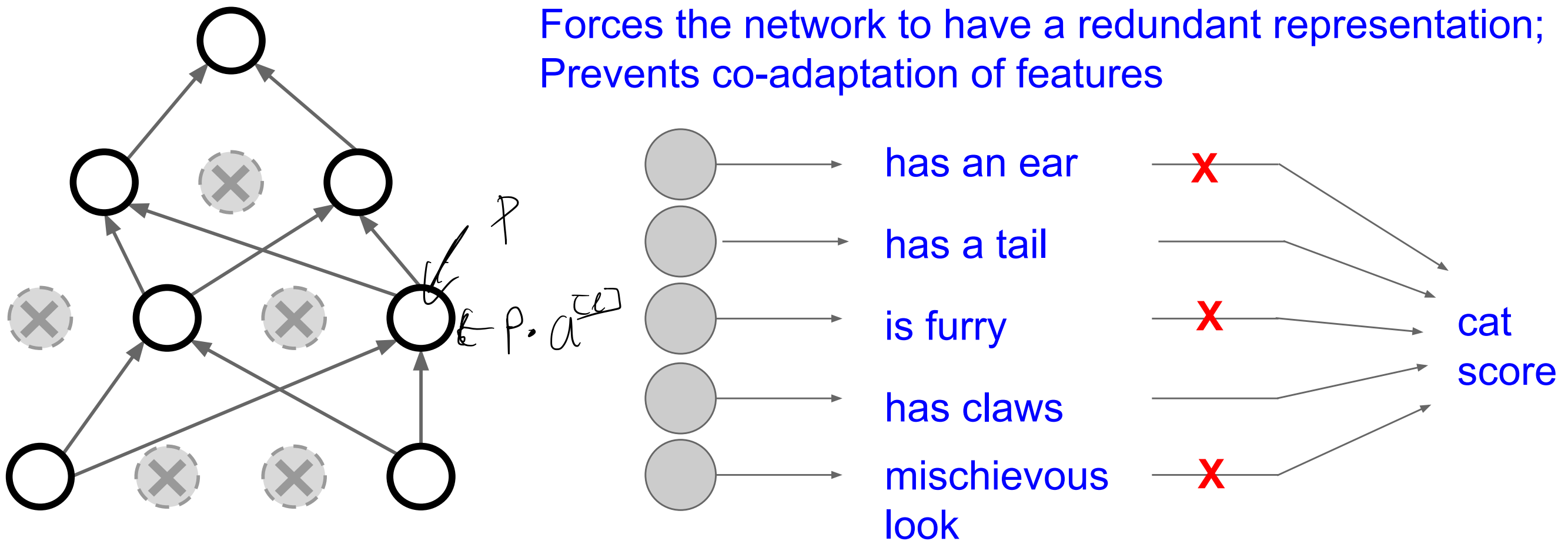
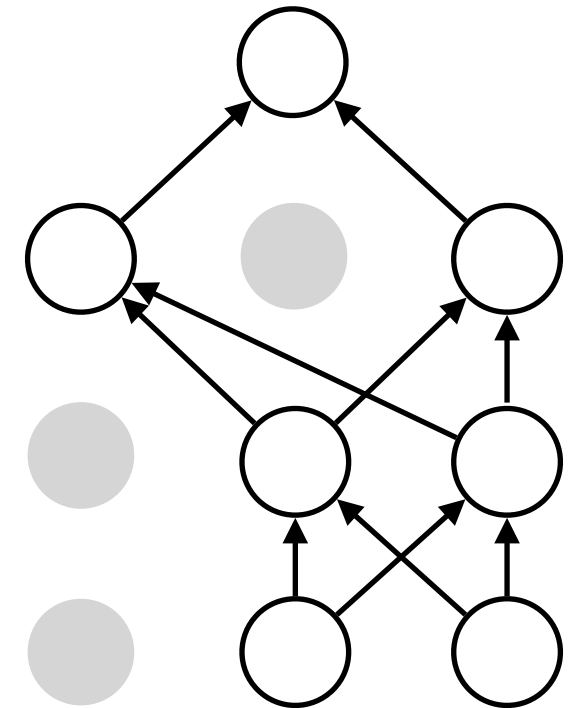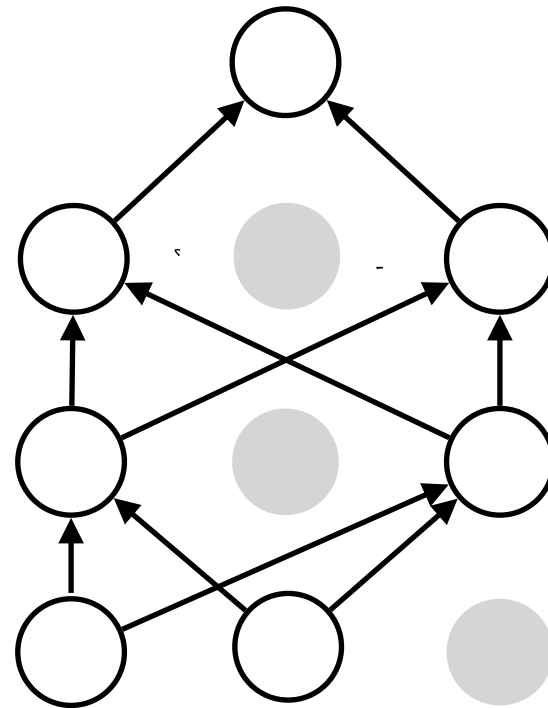# Dropout



- In each forward pass, randomly erase neurons

# Dropout
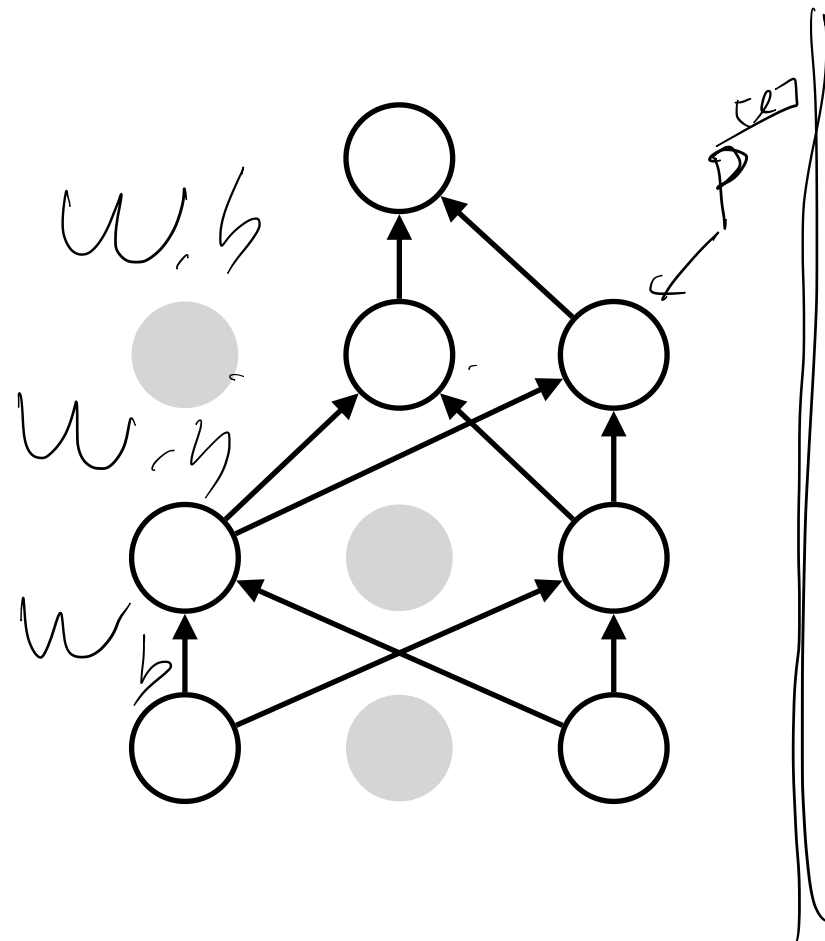
- Why can this be good?



Forces the network to have a redundant representation;
Prevents co-adaptation of features

has an ear   X

has a tail

is furry   X

has claws

mischievous look   X

cat score

# Ensemble of models

- Dropout is training a large ensemble of models (that share parameters).

- Each binary mask is one model

# Dropout: Test time

no dropout.
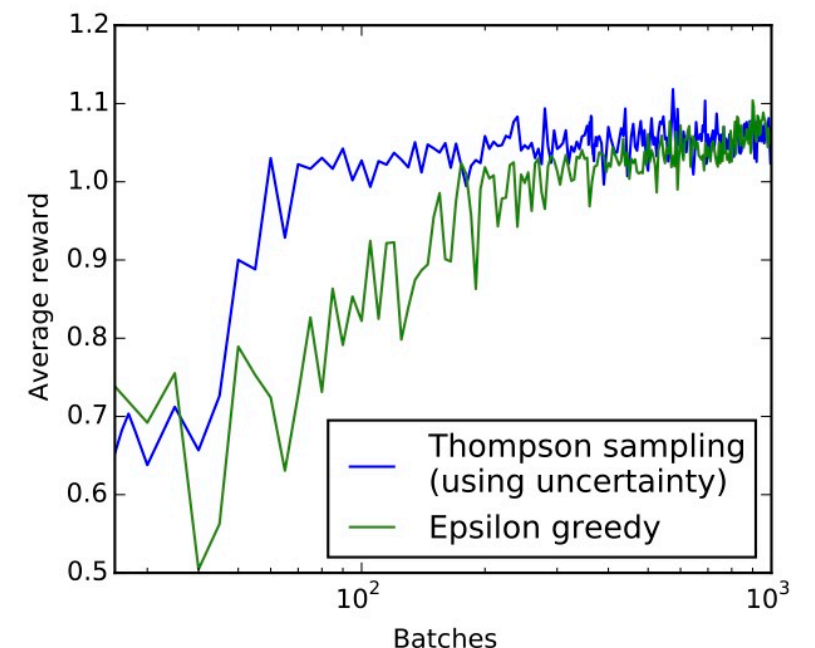
$$a^{[0]} = X$$
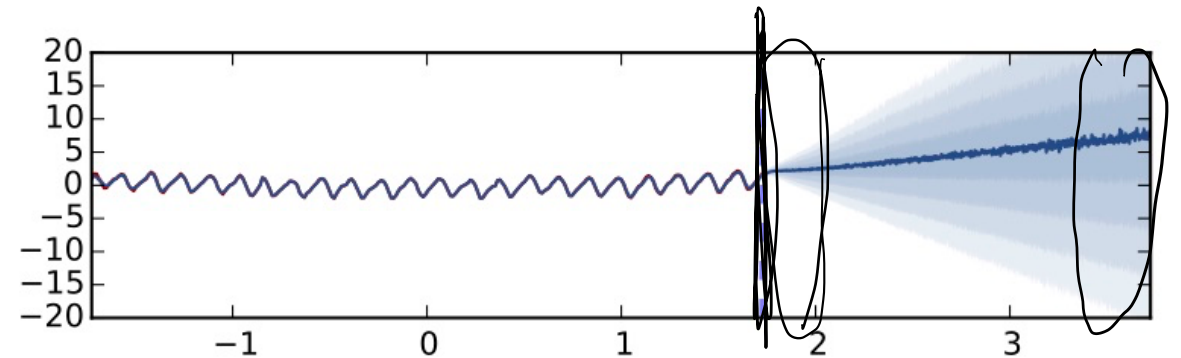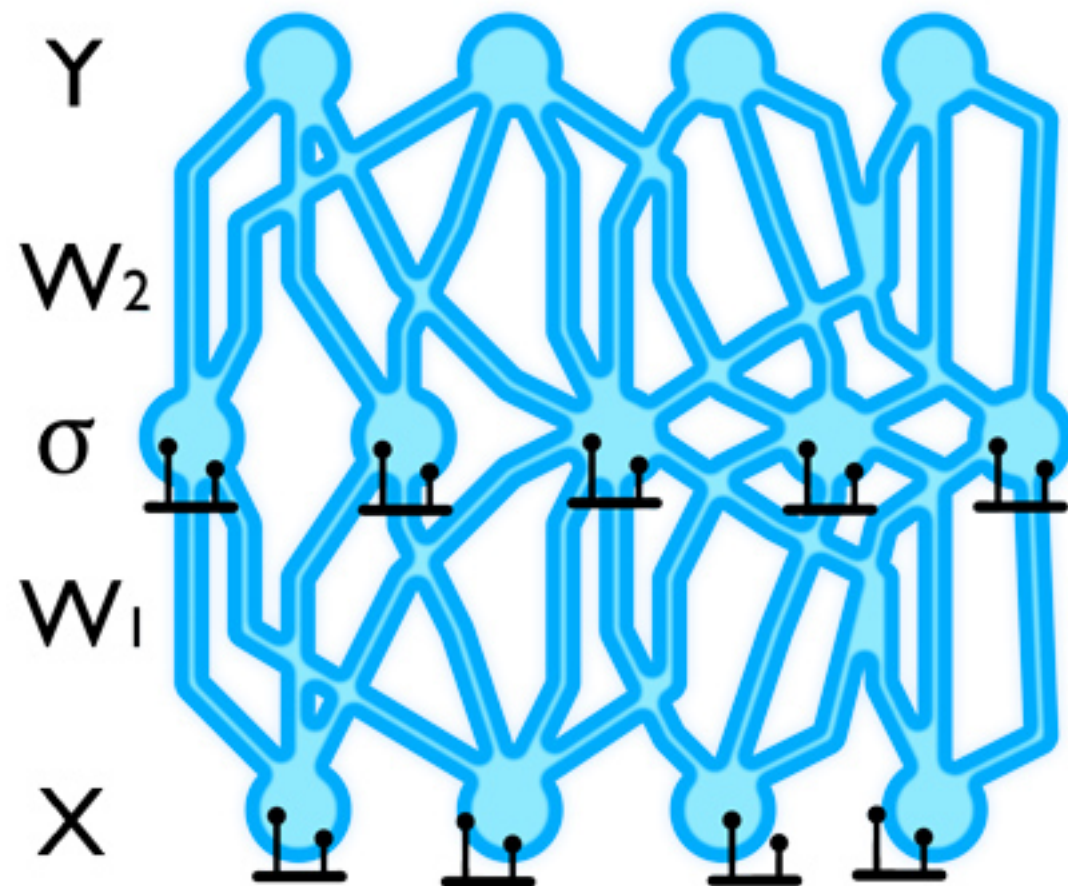
$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = \sigma^{[1]}(z^{[1]})$$

$$\Rightarrow \quad a^{[\ell]} = p^{[\ell]} \cdot \sigma^{[\ell]}(z^{[\ell]})$$

$$\hat{y} = a^{[L]}$$

- No dropout at test time

- scaling by dropout probability

# Dropout: uncertainty



- Dropout generates ensemble of models

- From the ensemble, estimate mean and variance of the output

Y. Gal and Z. Ghahrami, "Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning," https://arxiv.org/abs/1506.02142

10

# SGD and Early stopping

- SGD adds noises to the network -> a regularizer

- Early stopping

$$\mathcal{L}(\omega) = \sum_{i=1}^{n} \ell\left(Y^{(i)}, \hat{Y}^{(i)}\right)$$



Loss

Iteration

**Train**
**Val**

Accuracy

Stop training here

Iteration

# Batch Normalization

- Covariate shift

- "The change in the distributions of layers' inputs presents a problem because the layers need to continuously adapt to the new distribution. When the input distribution to a learning system changes, it is said to experience covariate shift"
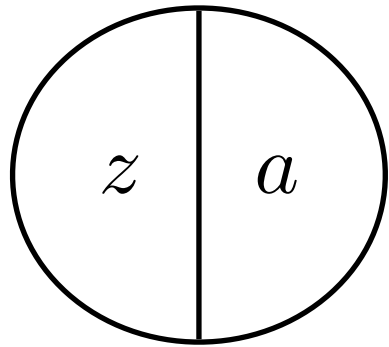
# Batch Normalization

$N(0, 1)$



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

$$a^{[\ell+1]} = \sigma(z^{[\ell+1]})$$

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

$$z_{norm}^{[\ell+1]} = \frac{z^{[\ell+1]} - \mu_z^{[\ell+1]}}{\sigma_z^{[\ell+1]}}$$

$$\tilde{z}^{[\ell+1]} = \gamma^{[\ell+1]} \odot z_{norm}^{[\ell+1]} + \beta^{[\ell+1]}$$

Vector. element-wise.

$$a^{[\ell+1]} = \sigma(\tilde{z}^{[\ell+1]})$$

13

# Training with Batch Norm.

- Original: $\{w^{[1]}, b^{[1]}, \ldots, w^{[L]}, b^{[L]}\}$

- With Batch Norm.: $\{w^{[1]}, b^{[1]}, \beta^{[1]}, \gamma^{[1]}, \ldots, w^{[L]}, b^{[L]}, \beta^{[L]}, \gamma^{[L]}\}$

$$w, b, \beta, \gamma \leftarrow w, b, \beta, \gamma - \eta \nabla \mathcal{L}(w, b, \beta, \gamma)$$

$$\underset{\tilde{\uparrow}}{\eta}$$

$$learning$$

$$\hat{\mu}^{[\ell]}, \hat{\sigma}^{[\ell]}$$

$$with \quad mini-Batch, \quad B = 8, 16, 32.$$

$$X^{\{1\}}, X^{\{2\}}, X^{\{3\}}, \ldots$$

$$X^{\{m\}} = \{ X^{\{m\}(1)}, \ldots, X^{\{m\}(B)} \}$$

$$\hat{\mu}^{[\ell]} = \frac{1}{B} \sum_{b=1}^{B} z^{[\ell]\{m\}(b)} \qquad \hat{\sigma}^{[\ell]} = \frac{1}{B-1} \sum_{b=1}^{B} (z^{[\ell]\{m\}(b)} - \hat{\mu}^{[\ell]})^2$$

14

# Batch Norm at test time

$$\mu^{[\ell]}, \sigma^{[\ell]}$$

$$\mu^{[\ell]} \approx \hat{\mu}^{[\ell]} \leftarrow \text{last mini-batch.}$$

$$\approx \hat{\hat{\mu}}^{[\ell]} \leftarrow \text{using all training data.}$$

$$\boxed{\hat{\mu}^{[\ell]}} = \frac{\eta^M \hat{\mu}^{[\ell]\{1\}} + \eta^{M-1} \hat{\mu}^{[\ell]\{2\}} + \cdots + \eta \hat{\mu}^{[\ell]\{M\}}}{\eta^M + \cdots \cdots + \eta} \qquad \eta < 1$$

$$\Rightarrow \boxed{\hat{\mu}^{[\ell]}} \Leftarrow (1-\eta)\boxed{\hat{\mu}^{[\ell]\{M\}}} + \eta \boxed{\hat{\mu}^{[\ell]}}$$

# Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.

- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.

- This has a slight regularization effect.