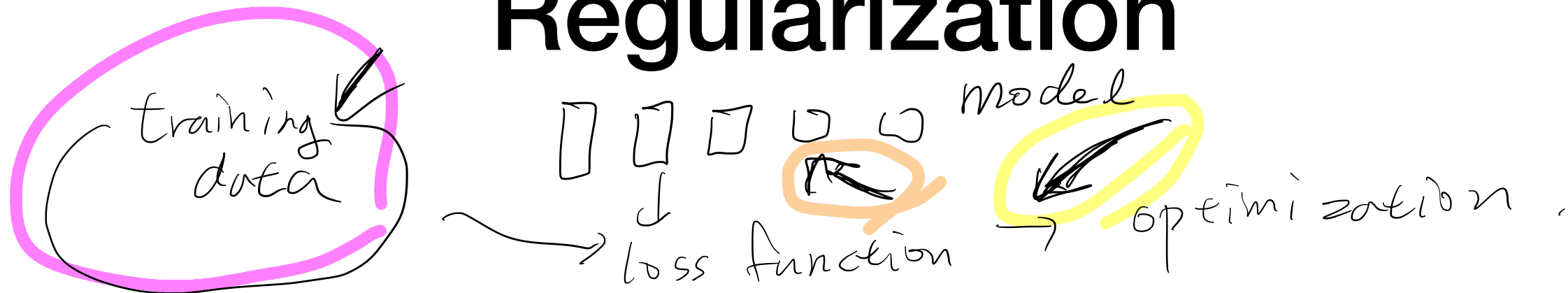


Regularization

Seyoung Yun

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

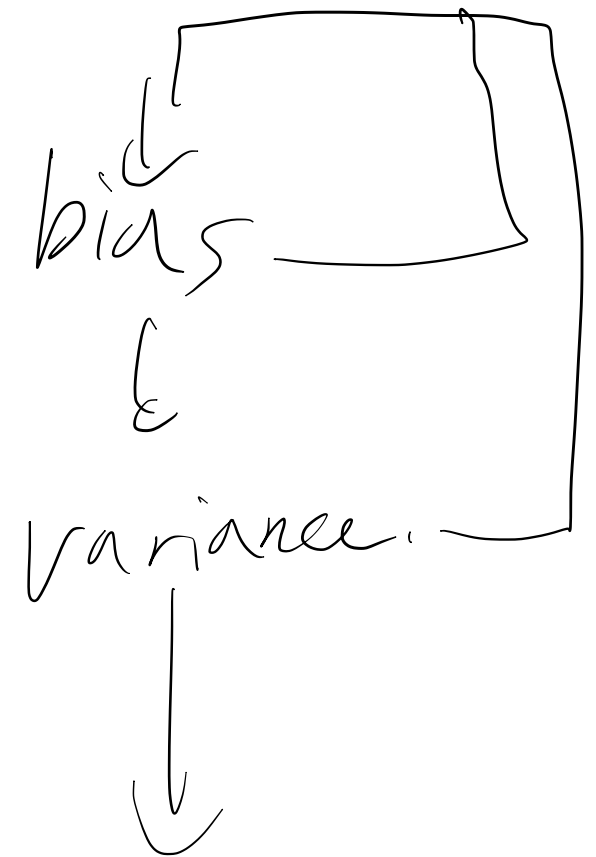
Regularization



- “A regularizer is anything that hurts the training process”

C. Zhung at ICLR2017 (<https://www.youtube.com/watch?v=kCj51pTQPKI>)

- data augmentation
- weight decay - with an additional cost
- dropout - by adding random noise



Linear Regression

- RSS: cost of linear regression

$$\mathcal{L}(w, b) = \sum_{i=1}^m (\underbrace{y^{(i)}}_{\text{target}} - \underbrace{w^\top x^{(i)} + b}_{\text{linear function}})^2$$

- 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

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{i=1}^L \|w^{[i]}\|_F^2$$

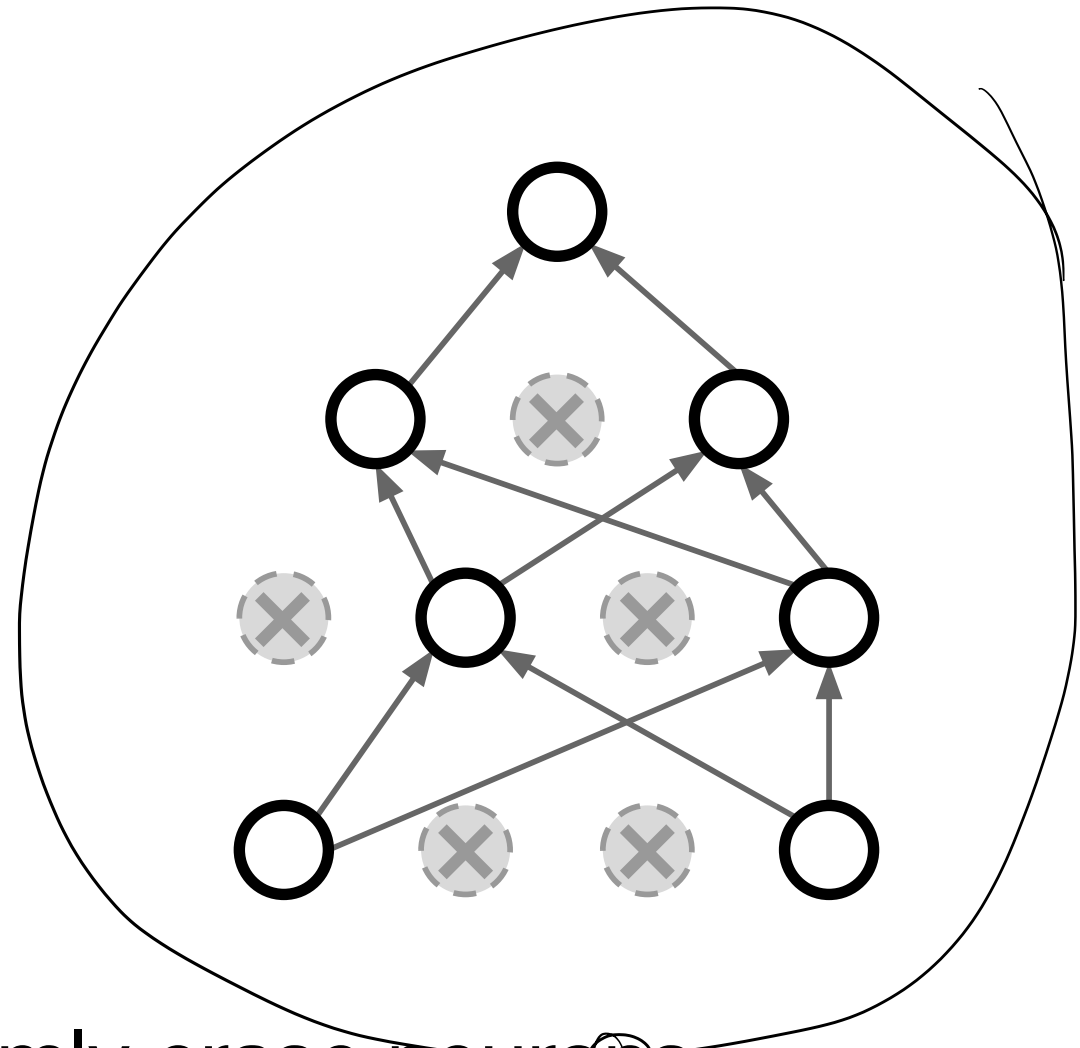
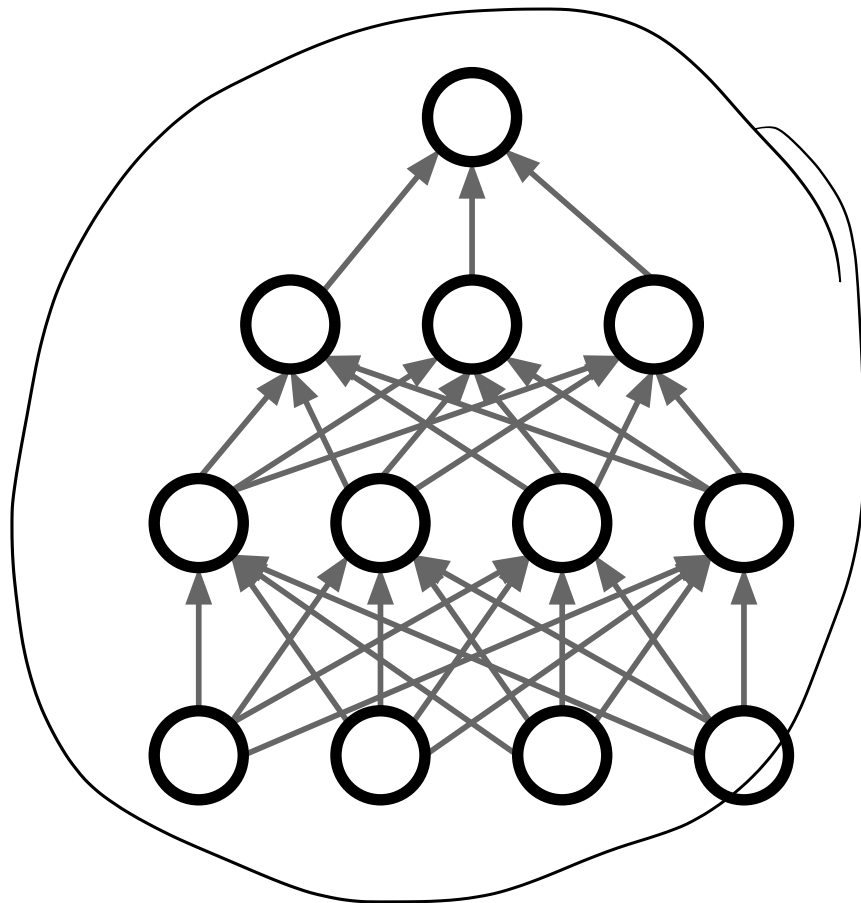
$$\|w^{[L]}\|_F^2 = \sum_j \sum_k (w_{jk}^{[L]})^2$$

$$\Rightarrow \nabla_{w^{[L]}} J = \frac{1}{m} \sum_{i=1}^m \nabla_{w^{[L]}} \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{m} \cdot w^{[L]}$$

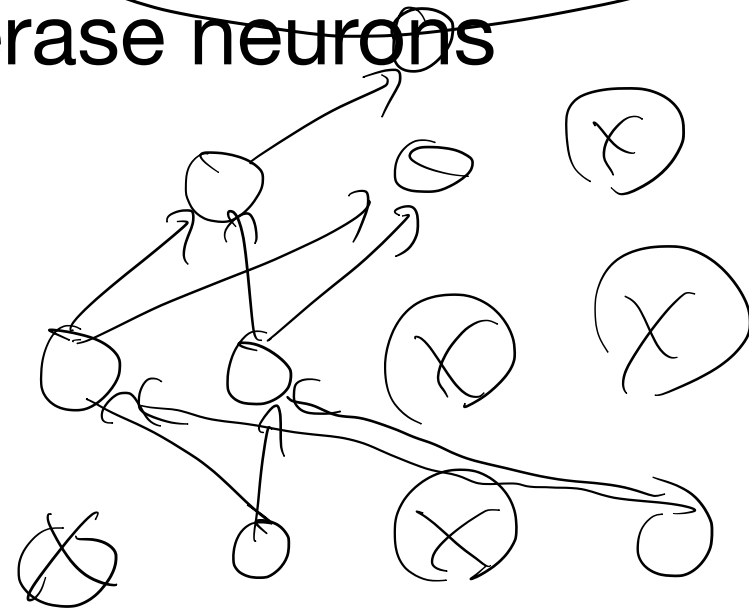
Back prop.

$$\begin{aligned} w^{[L]}(t+1) &= w^{[L]}(t) - \gamma \cdot \nabla_{w^{[L]}} J(w^{[L]}(t)) \\ &= \left(1 - \frac{\gamma \lambda}{m}\right) w^{[L]}(t) - \underbrace{\frac{\gamma}{m} \sum_{i=1}^m \nabla_{w^{[L]}} \ell(\hat{y}^{(i)}, y^{(i)})}_{\text{Back prop.}} \end{aligned}$$

Dropout

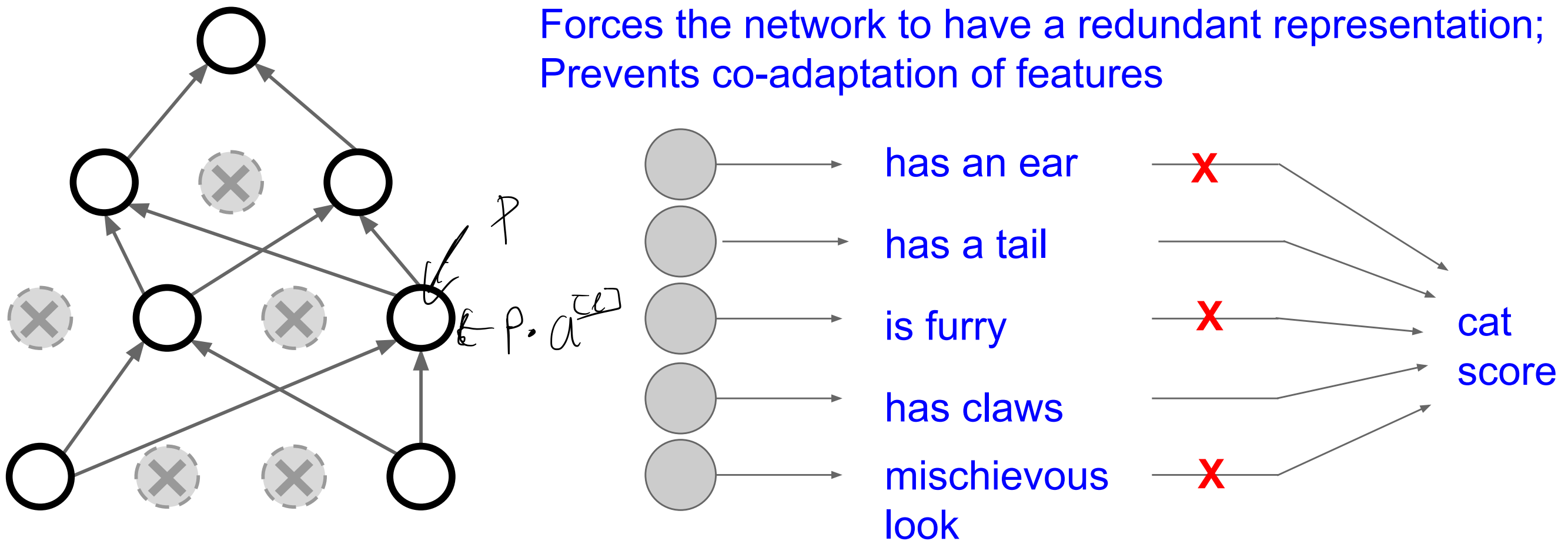


- In each forward pass, randomly erase neurons

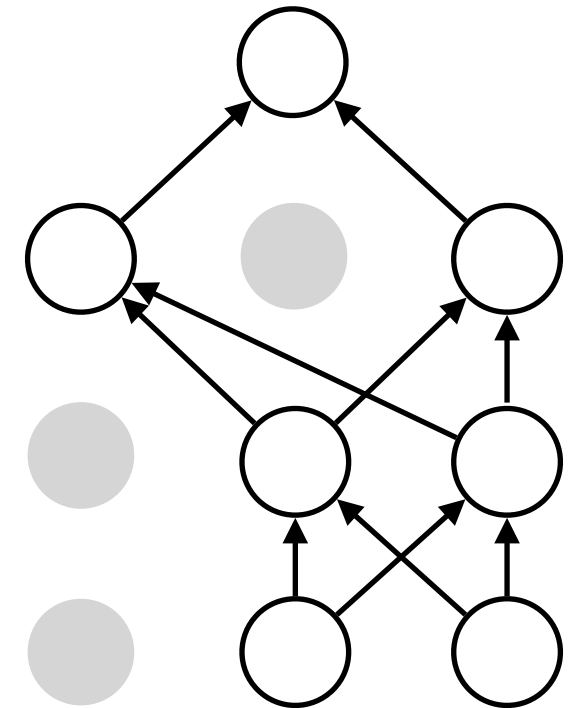
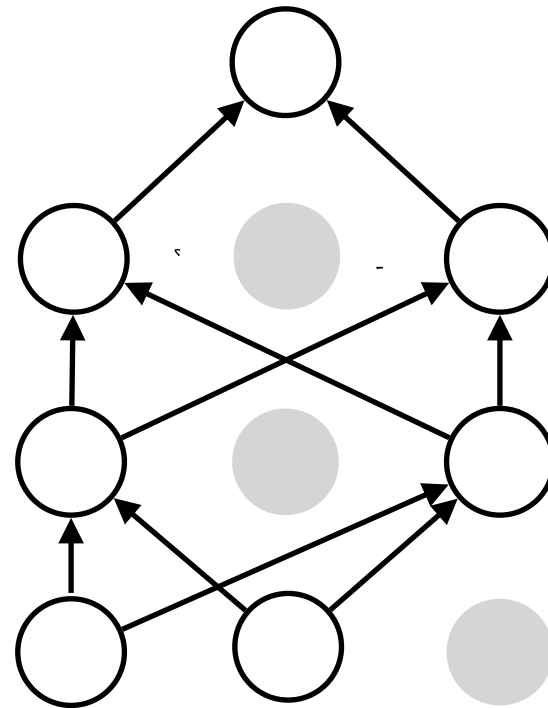
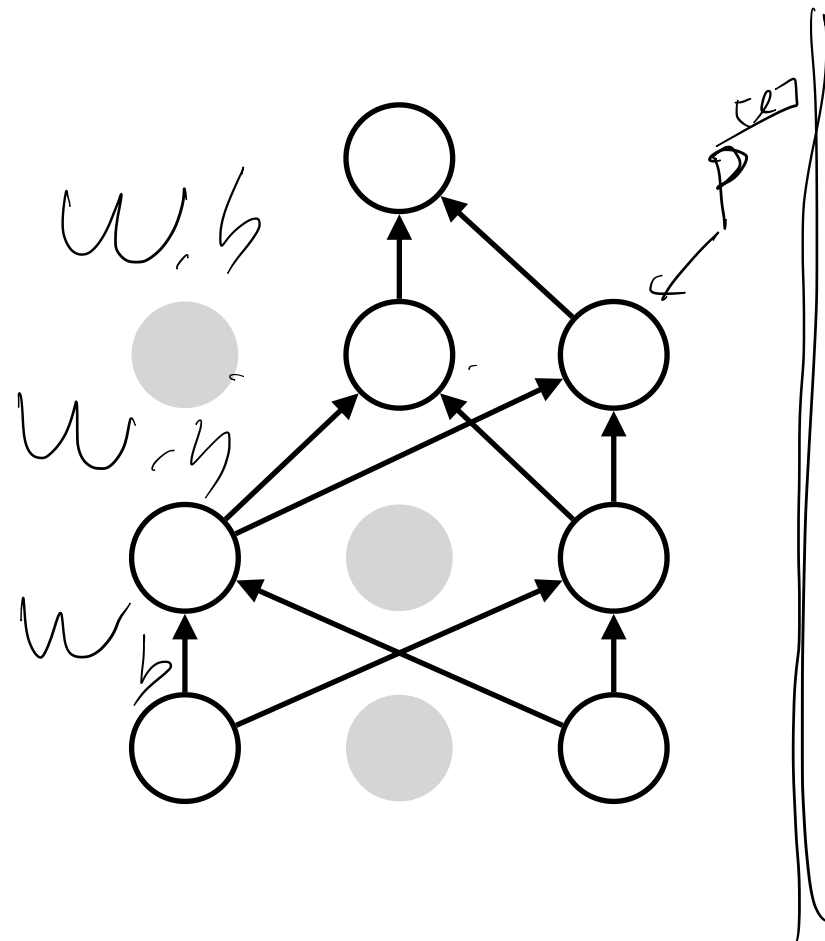


Dropout

- Why can this be good?



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 -

$$a^{[0]} = X$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

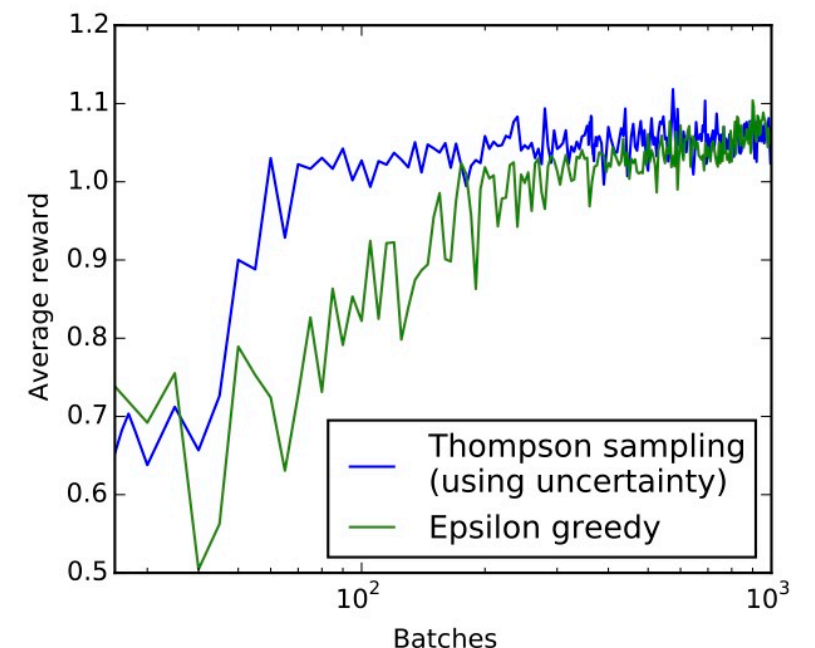
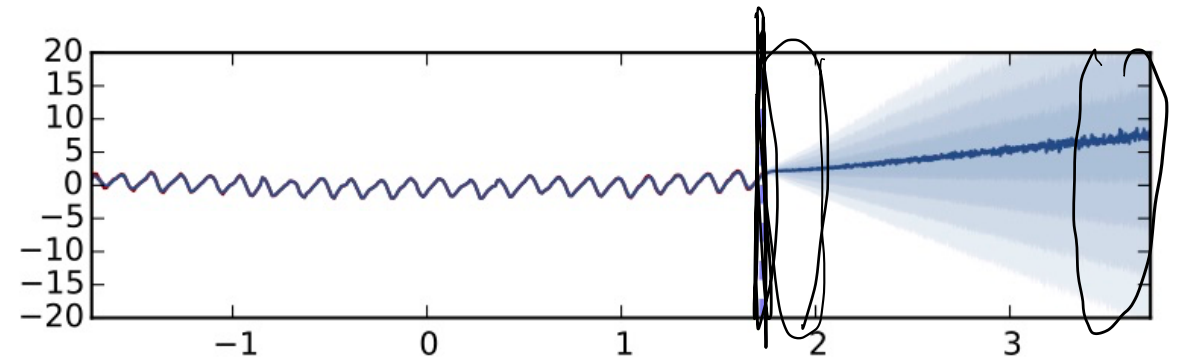
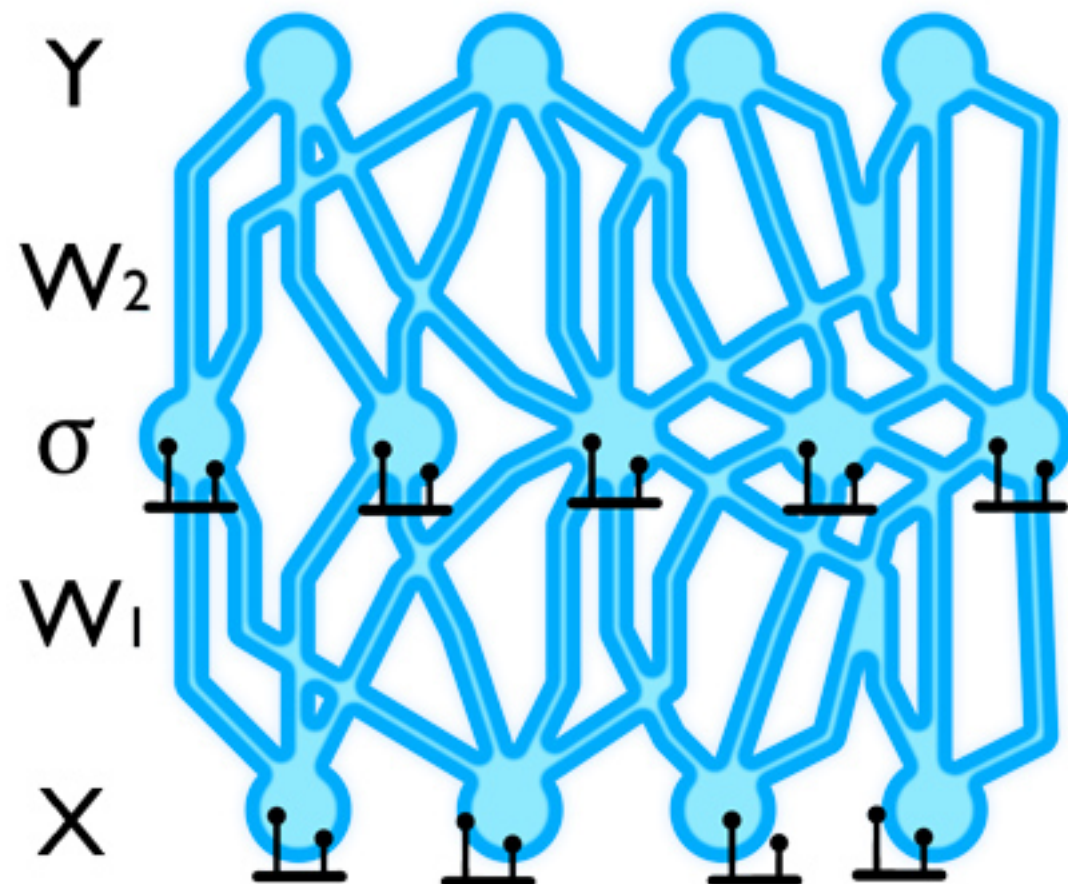
$$\hat{y} = a^{[L]}$$

\Rightarrow

$$a^{[L]} = p^{[L]} \cdot \sigma^{[L]}(z^{[L]})$$

- No dropout at test time
- scaling by dropout probability

Dropout: uncertainty



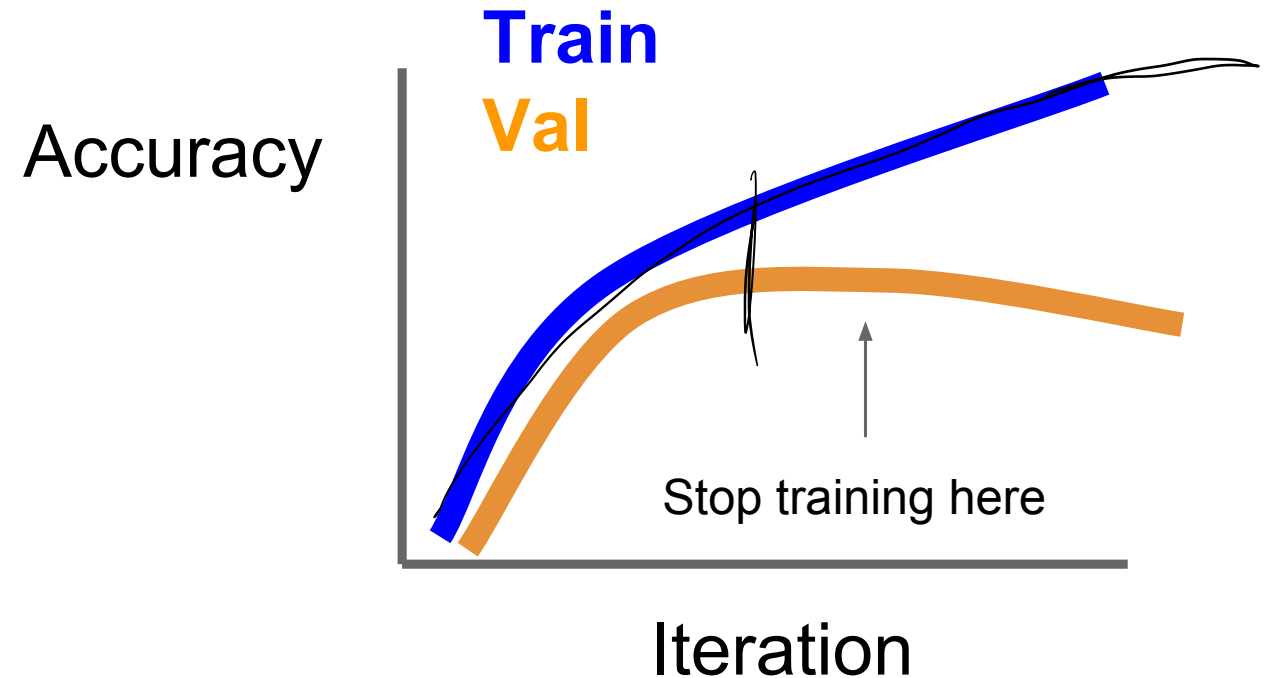
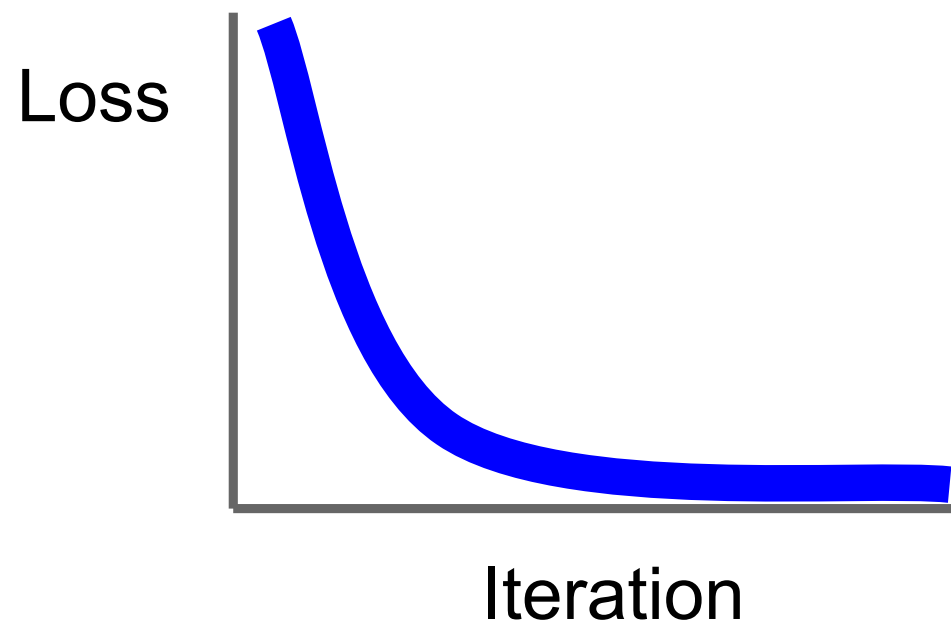
- Dropout generates ensemble of models
- From the ensemble, estimate mean and variance of the output

SGD and Early stopping

- SGD adds noises to the network -> a regularizer

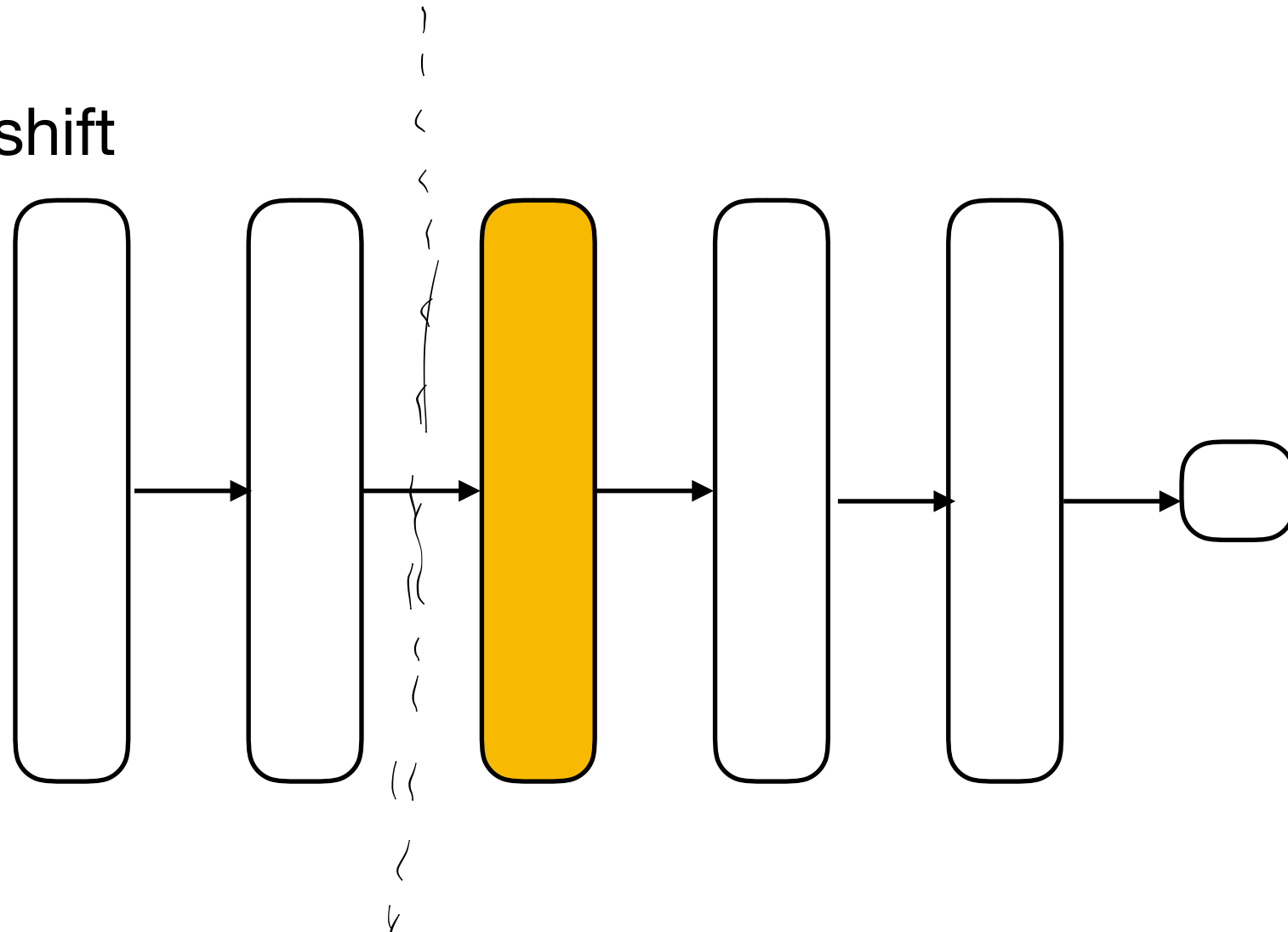
- Early stopping

$$\mathcal{L}(w) = \sum_{\tilde{i}=1}^n \ell(y^{(\tilde{i})}, \hat{y}^{(\tilde{i})})$$



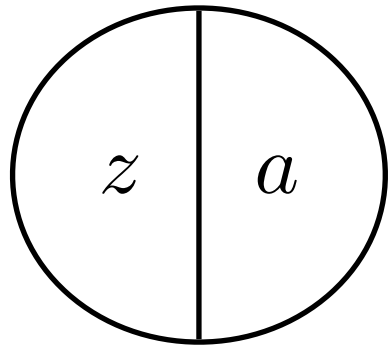
Batch Normalization

- Covariate shift



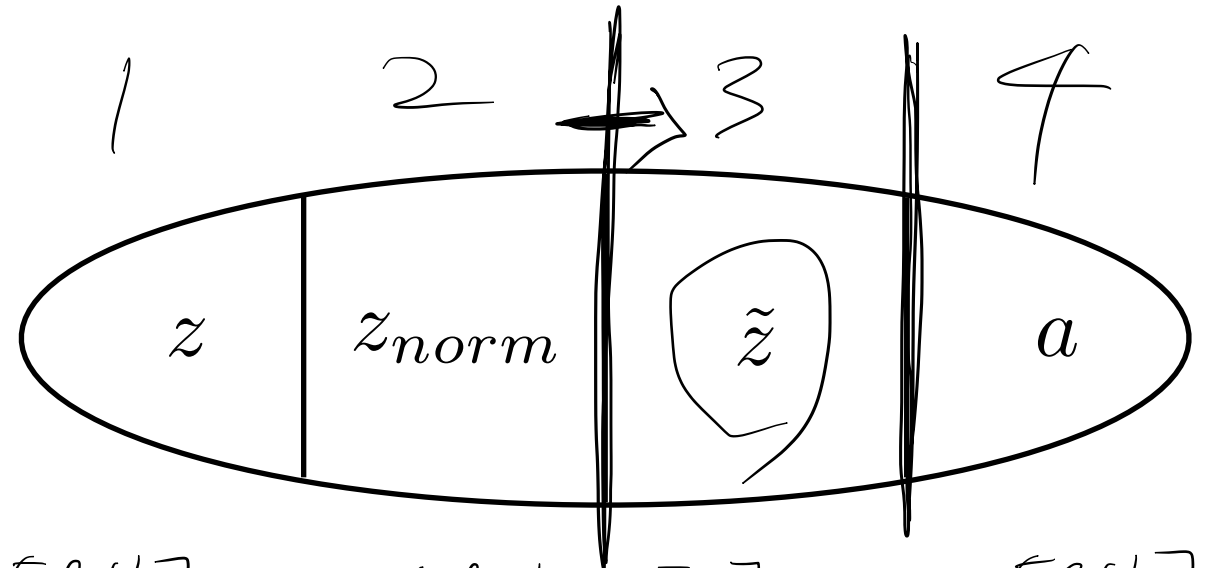
- “The change in the distributions of layers’ inputs presents a problem because the layers need to continuously adapt to the new distribution. When the input distribution to a learning system changes, it is said to experience covariate shift”

Batch Normalization



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

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



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

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

$$\tilde{z}^{[l+1]} = \gamma^{[l+1]} \otimes z_{\text{norm}}^{[l+1]} + \beta^{[l+1]}$$

\uparrow Vector. element-wise

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

Training with Batch Norm.

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

Batch Norm at test time

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.