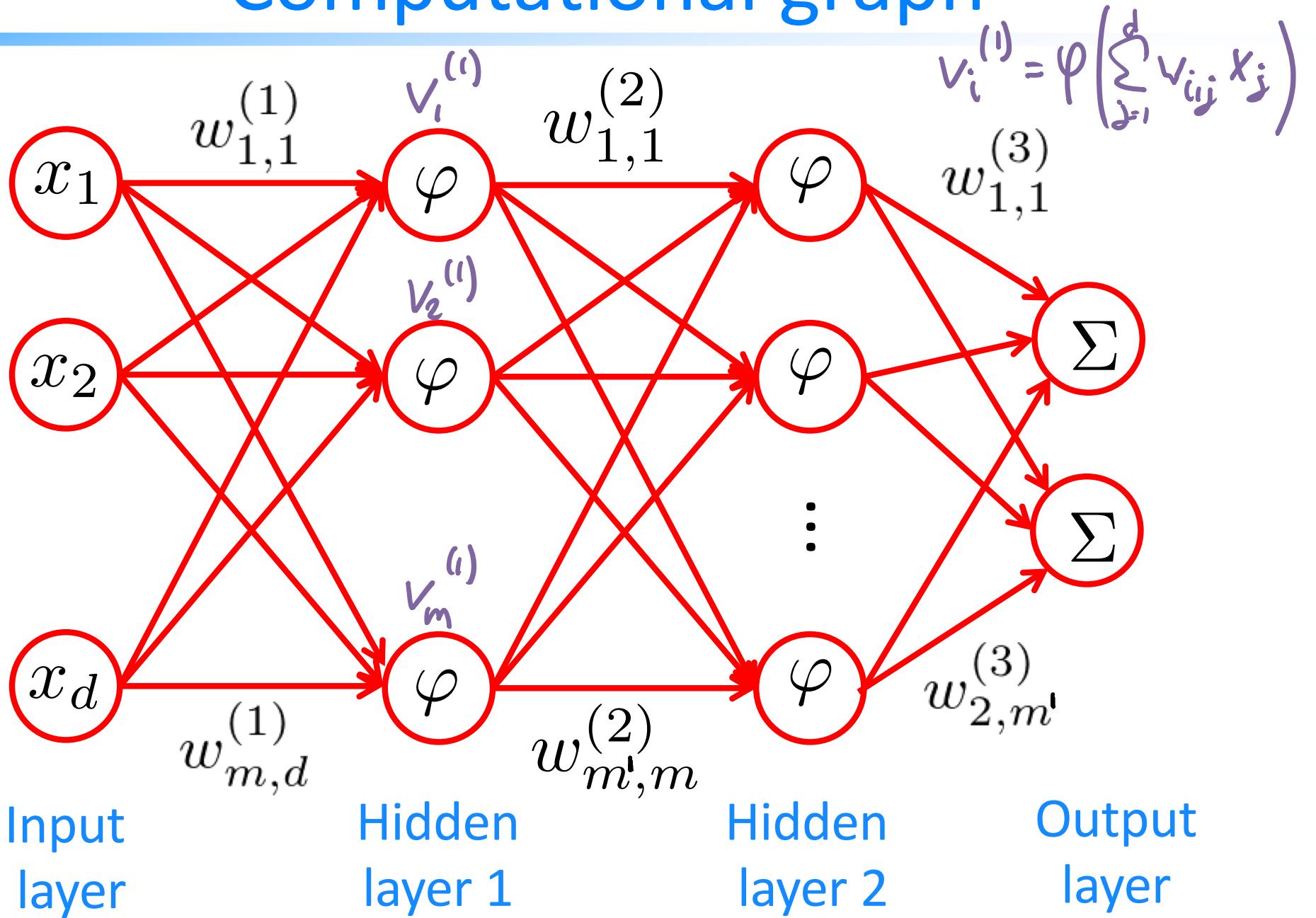


Introduction to Machine Learning

Neural networks (continued)

Prof. Andreas Krause
Learning and Adaptive Systems (las.ethz.ch)

Computational graph



How can we train the weights?

- Given data set $D = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$
want to optimize weights $\mathbf{W} = (\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)})$
- How do we measure and optimize goodness of fit?
→ Apply **loss function** (e.g., Perceptron loss, multi-class hinge loss, square loss, etc.) to output

$$\ell(\mathbf{W}; \mathbf{y}, \mathbf{x}) = \ell(\mathbf{y} - f(\mathbf{x}, \mathbf{W}))$$

→ Then **optimize the weights** to minimize loss over D

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^n \ell(\mathbf{W}; \mathbf{x}_i, y_i)$$

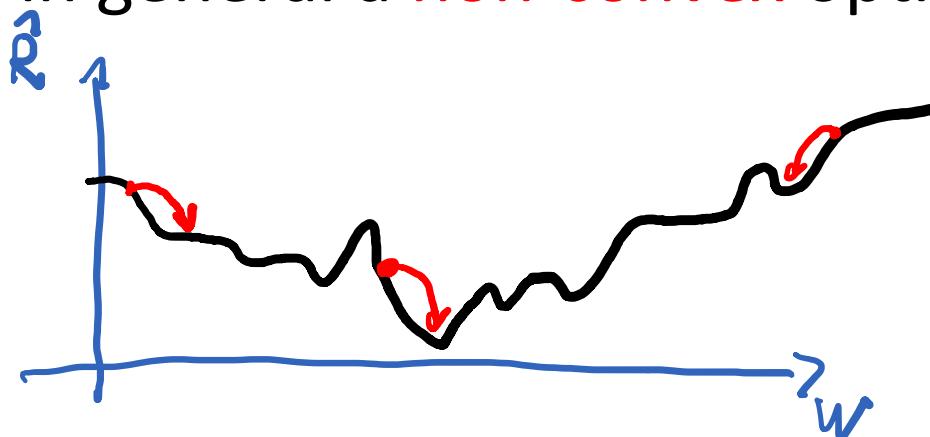
How do we optimize over weights?

- Want to do Empirical Risk Minimization

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^n \ell(\mathbf{W}; \mathbf{x}_i, y_i)$$

R(w)

- I.e., jointly **optimize over all weights for all layers** to minimize loss over the training data
- This is in general a **non-convex** optimization problem



- Nevertheless, can try to find a local optimum

Stochastic gradient descent for ANNs

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^n \ell(\mathbf{W}; \mathbf{x}_i, y_i)$$

- Initialize weights \mathbf{W}
- For $t = 1, 2, \dots$
 - Pick data point $(\mathbf{x}, \mathbf{y}) \in D$ uniformly at random
 - Take step in negative gradient direction

$$\mathbf{W} \leftarrow \mathbf{W} - \eta_t \nabla_{\mathbf{W}} \ell(\mathbf{W}; \mathbf{y}, \mathbf{x})$$

Backpropagation: Matrix form

- For the output layer

- Compute „error“ $\delta^{(L)} = \mathbf{l}'(\mathbf{f}) = [l'(f_1), \dots, l'(f_p)]$

- Gradient: $\nabla_{\mathbf{W}^{(L)}} \ell(\mathbf{W}; \mathbf{y}, \mathbf{x}) = \delta^{(L)} \mathbf{v}^{(L-1)T}$

- For each hidden layer $\ell = L - 1 : -1 : 1$

- Compute „error“ $\delta^{(\ell)} = \varphi'(\mathbf{z}^{(\ell)}) \odot (\mathbf{W}^{(\ell+1)T} \delta^{(\ell+1)})$

pointwise multiplication

- Gradient $\nabla_{\mathbf{W}^{(\ell)}} \ell(\mathbf{W}; \mathbf{y}, \mathbf{x}) = \delta^{(\ell)} \mathbf{v}^{(\ell-1)T}$

Some questions

- How can we compute the gradient?
- How should we initialize the weights?
- When should we terminate?
- How do we choose parameters (number of units/layers/activation functions/learning rate/ ...)?
- What about overfitting?

Initializing weights

- Since optimization problem is non-convex, initialization matters!
- What happens if we pick inappropriate weights?

Propagation of variance (ReLU)

$$v_i^{(l-1)} \rightarrow v_i^l$$

⋮

$$v_{n_m}^{(l-1)} \rightarrow v_i^l$$

$$z_i^{(l)} = \sum_{j=1}^{n_m} w_{ij} v_j^{(l-1)} \Rightarrow v_i^{(l)} = \varphi(z_i^{(l)})$$

$\varphi(z) = \max(0, z)$

Assume: $E[x_i] = E[v_i^{(0)}] = 0, \text{Var}[x_i] = \text{Var}[v_i^{(0)}]$

Further: x_1, \dots, x_d independent.

Also: $w_{ij} \sim N(0, \sigma^2)$ independent

Sps x, y indep., - mean 0

- (a) if y mean 0: $\text{Var}(k y) = k \text{Var}(y) | k \neq 0$
- (b) if y mean $\neq 0$: $\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y)$

For ReLU and symmetric y :

~~$E[\varphi(y)^2] = \frac{1}{2} \text{Var}(y)$~~

$$E[z_i^{(l)}] = E\left[\sum_{j=1}^{n_m} w_{ij} x_j\right] = \sum_j E[w_{ij} x_j] = \sum_j E[w_{ij}] E[x_j] = 0$$

$$\text{Var}[z_i^{(l)}] = \sum_j \text{Var}[w_{ij} x_j] \stackrel{(a)}{=} \sum_j \text{Var}(w_{ij}) \text{Var}(x_j) = n \cdot \sigma^2 \quad \text{First layer}$$

$$\text{Var}[z_i^{(l)}] = \sum_j \text{Var}[w_{ij} v_j^{(l-1)}] = \sum_j \underbrace{\text{Var}(w_{ij})}_{\sigma^2} \underbrace{\frac{1}{2} \text{Var}(z_j^{(l-1)})}_{\frac{1}{2} \text{Var}(z_i^{(l-1)})} = \boxed{\frac{1}{2} n_m \sigma^2 := 1} \Rightarrow \sigma^2 = \frac{2}{n_m}$$

Initializing weights

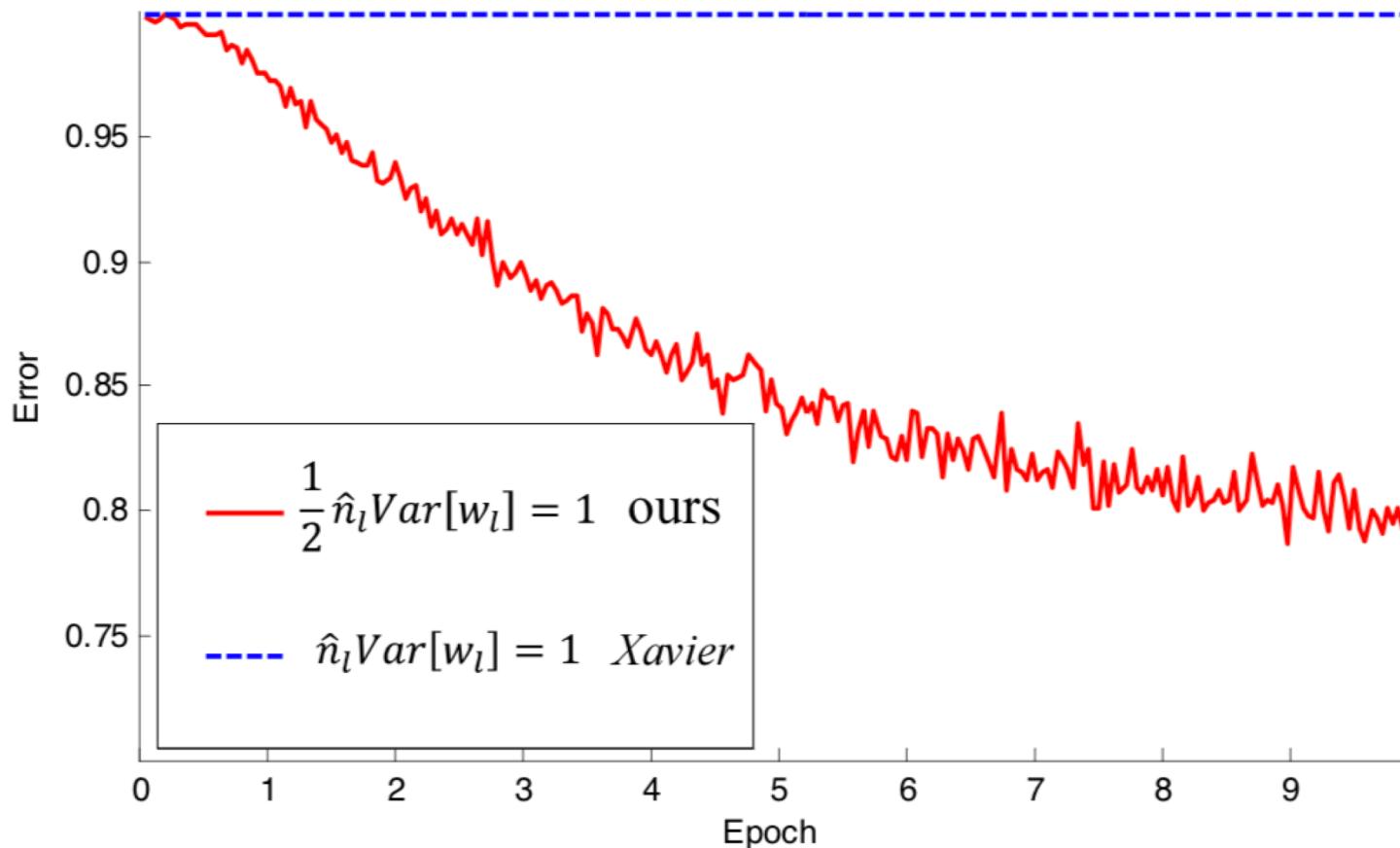
- Goal: Keep variance of weights approximately constant across layers to avoid vanishing and exploding gradients
- Usually random initialization works well, e.g.,

$$\text{Glorot (tanh): } w_{i,j} \sim \mathcal{N}(0, 1/(n_{in}))$$

$$w_{i,j} \sim \mathcal{N}(0, 2/(n_{in} + n_{out}))$$

$$\text{He (ReLU): } w_{i,j} \sim \mathcal{N}(0, \underbrace{2/n_{in}}_{\text{green underline}})$$

Effects of (im)proper initialization



[He, Zhang, Ren, Sun 2015]

Learning rate

- To implement SGD rule

$$\mathbf{W} \leftarrow \mathbf{W} - \eta_t \nabla_{\mathbf{W}} \ell(\mathbf{W}; \mathbf{y}, \mathbf{x})$$

need to choose learning rate η_t

- Usually good idea to start with fixed (small) learning rate, and decrease slowly after some iterations, e.g.,

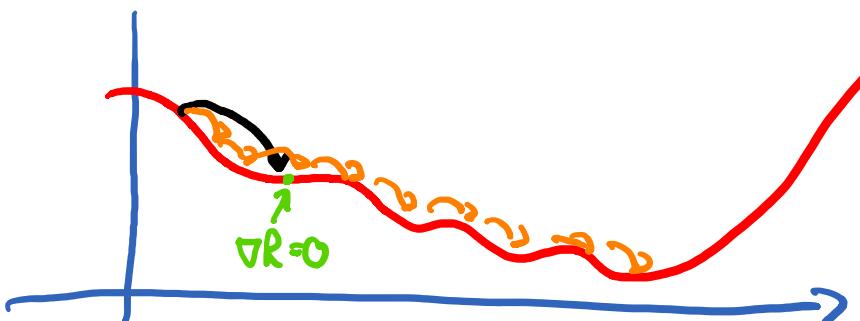
$$\eta_t = \min(0.1, 100/t)$$

- Often use piecewise constant learning rate „schedules“ (i.e., drop after some # of epochs)



Learning with momentum

- Common extension to training with SGD
- Can help to escape local minima
- *Idea:* Move not only into direction of gradient, but also in direction of last weight update

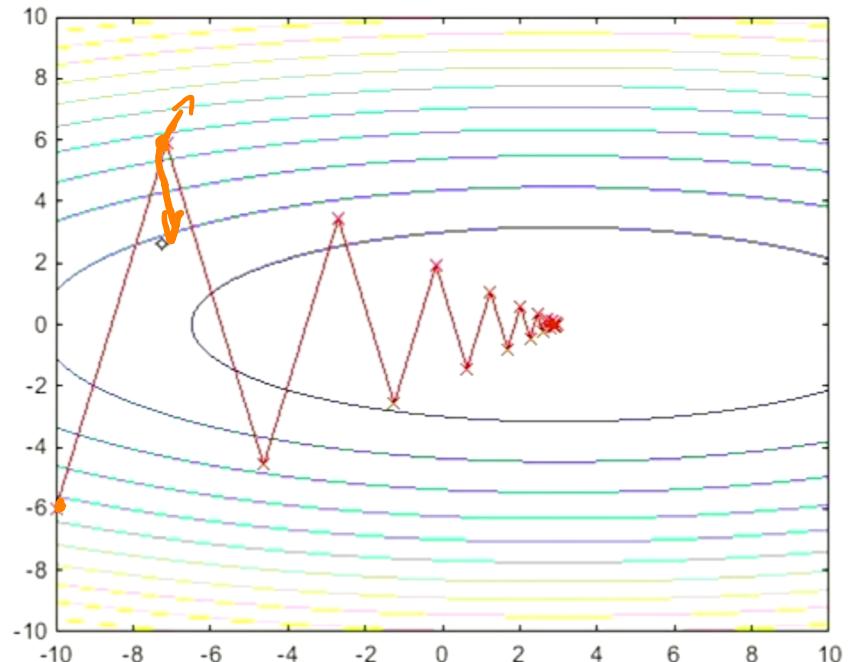
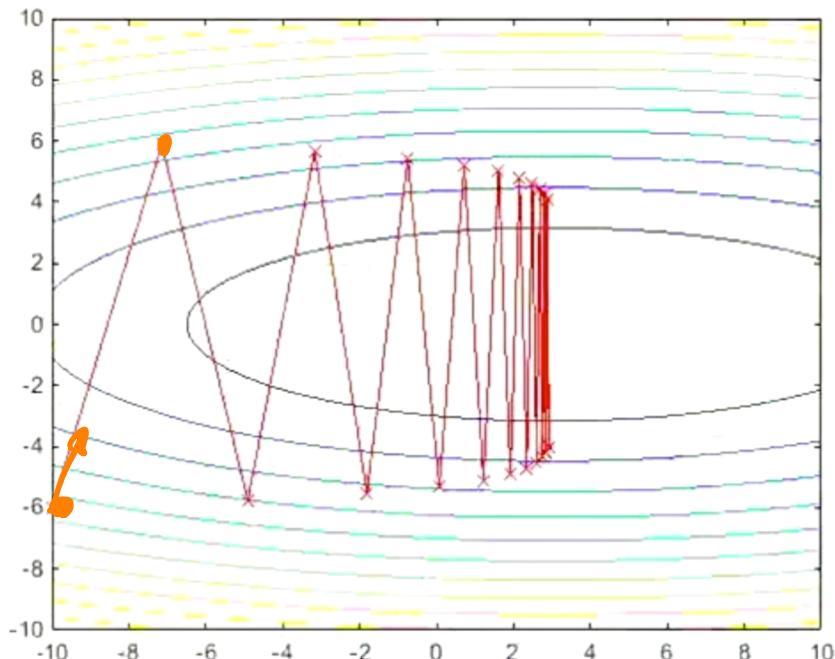


- Updates:
$$a \leftarrow m \cdot a + \eta_t \nabla_{\mathbf{W}} \ell(\mathbf{W}; \mathbf{y}, \mathbf{x})$$

$$\mathbf{W} \leftarrow \mathbf{W} - a$$

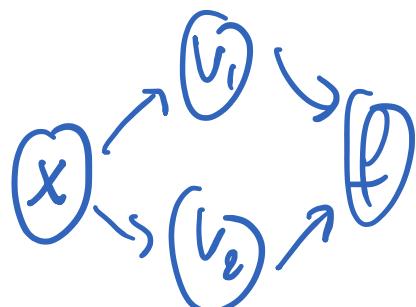
Learning with momentum

- Preventing oscillation



Weight-space symmetries

- Multiple distinct weights compute the same predictions



$$v_i = \varphi(w_i; x)$$

$$f = w_1' v_1 + w_2' v_2$$

$$w = [w_1, w_1'] ; , \bar{w} = [\bar{w}_1, \bar{w}_1'] ;$$

Sps. $\bar{w}_2 = w_1$, $\bar{w}_1 = w_2$, $\bar{w}_1' = w_2'$, $\bar{w}_2' = w_1'$

Then: $f(x; \bar{w}) = f(x; w)$

trah activations.

$$\varphi(z) = -\varphi(-z)$$

- Multiple local-minima can be equivalent in terms of input-output mapping $f(x; \bar{w}) = f(x; w)$

Avoiding overfitting

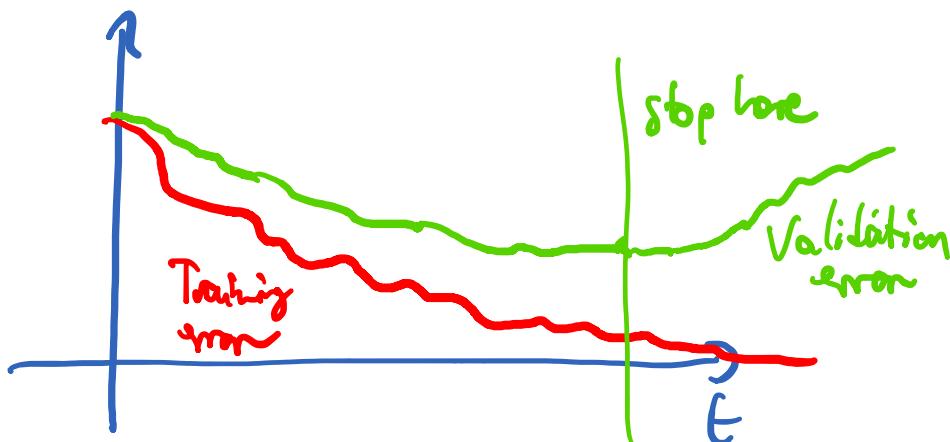
- Neural networks have many parameters
→ Potential danger of overfitting
- Countermeasures
 - Early stopping: Don't run SGD until convergence
 - Regularization: Add penalty term to keep weights small

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^n \ell(\mathbf{W}; \mathbf{x}_i, y_i) + \lambda ||W||_F^2$$

- „Dropout“

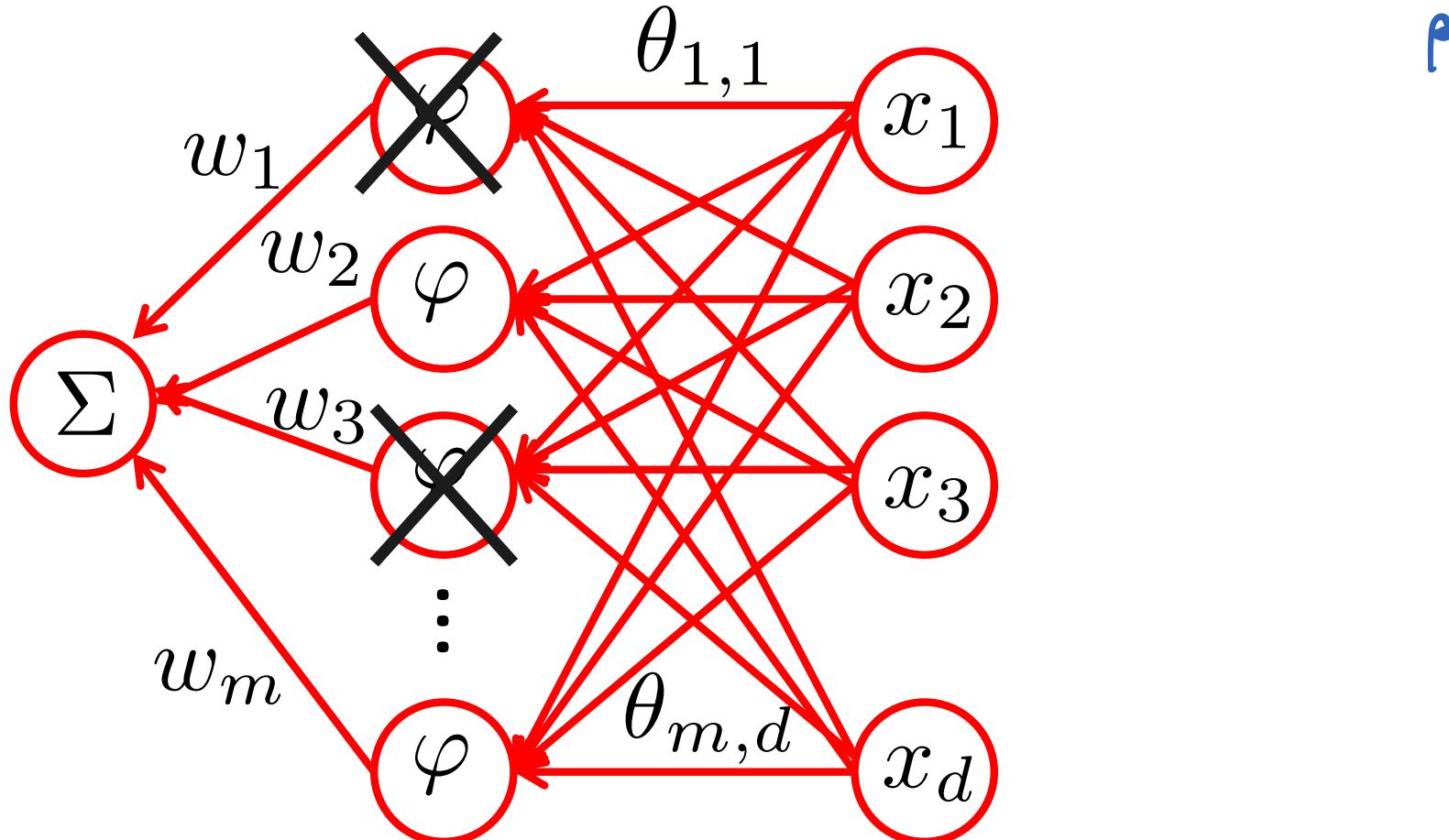
Early stopping

- In general, might not want to run training until weights converge → Overfitting!
- One possibility:
 - Monitor prediction performance on validation set;
 - stop training once validation error starts to increase



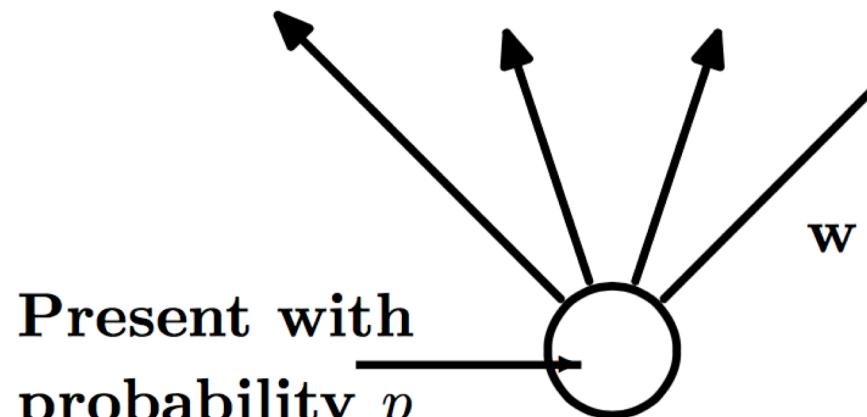
Dropout regularization

- Key idea: randomly ignore („drop out“) hidden units during each iteration of SGD with probability $1/2$

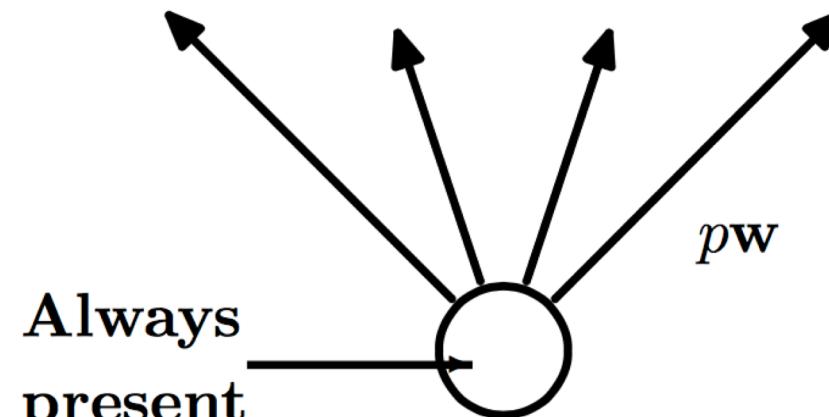


- After training, half the weights to compensate

Dropout: Train vs Test



(a) At training time

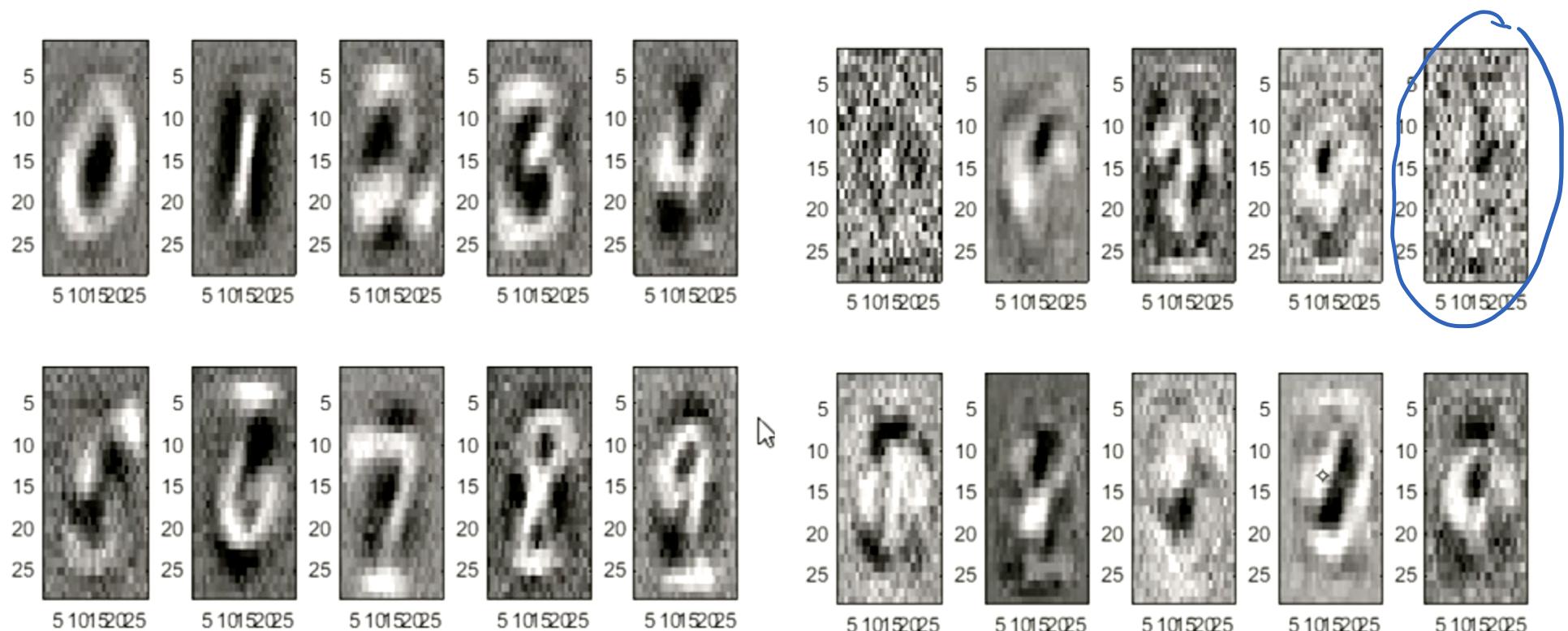


(b) At test time

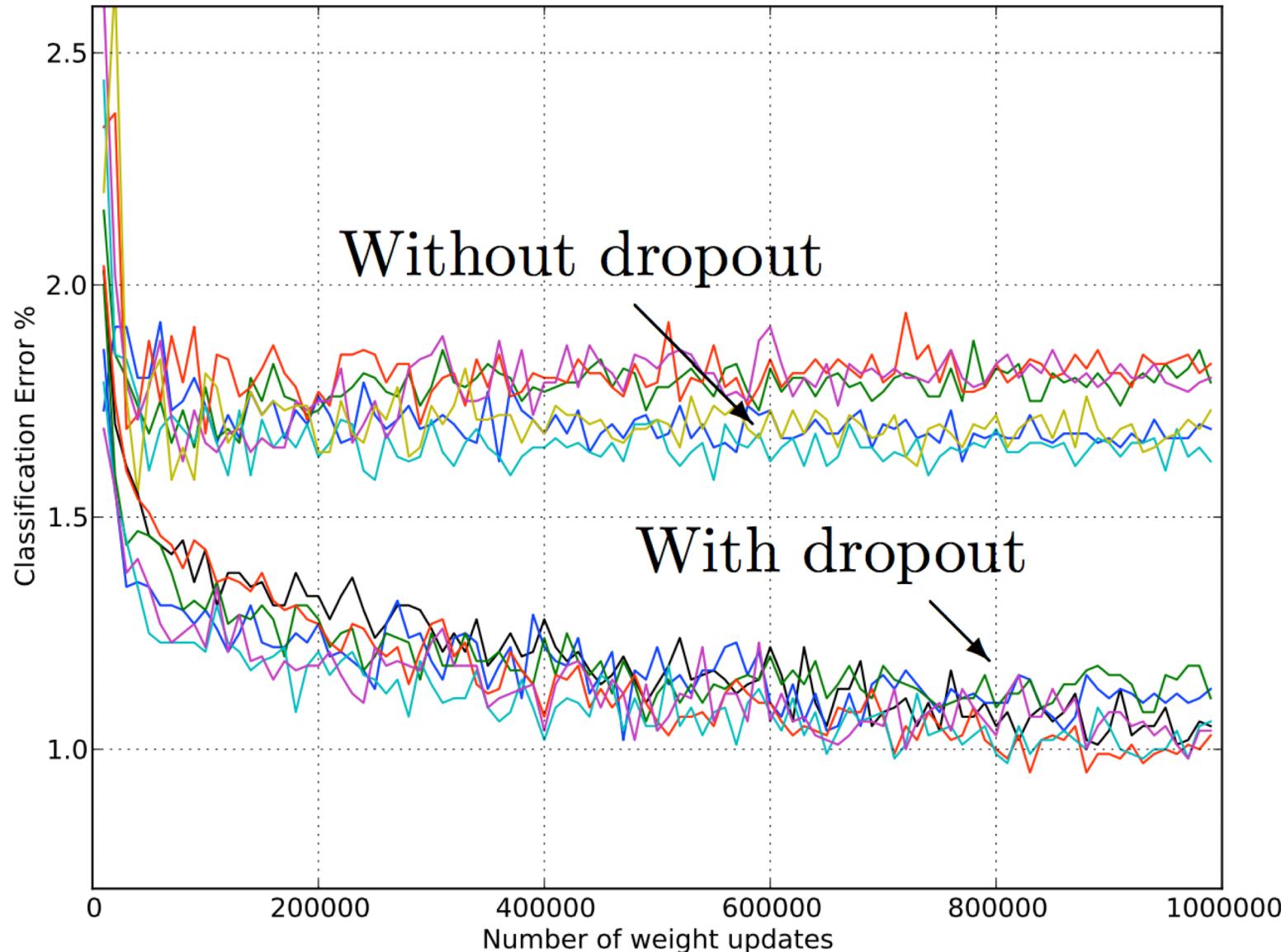
$$\begin{aligned} E[\varepsilon] &= E[w^T x] \\ &= E\left[\sum_i w_i x_i \mid i \text{ is selected}\right] \\ &\approx \sum_i w_i x_i p = p \cdot w^T x = (\rho w)^T x \end{aligned}$$

[Srivastava et al JMLR '14]

Ovefitting illustration

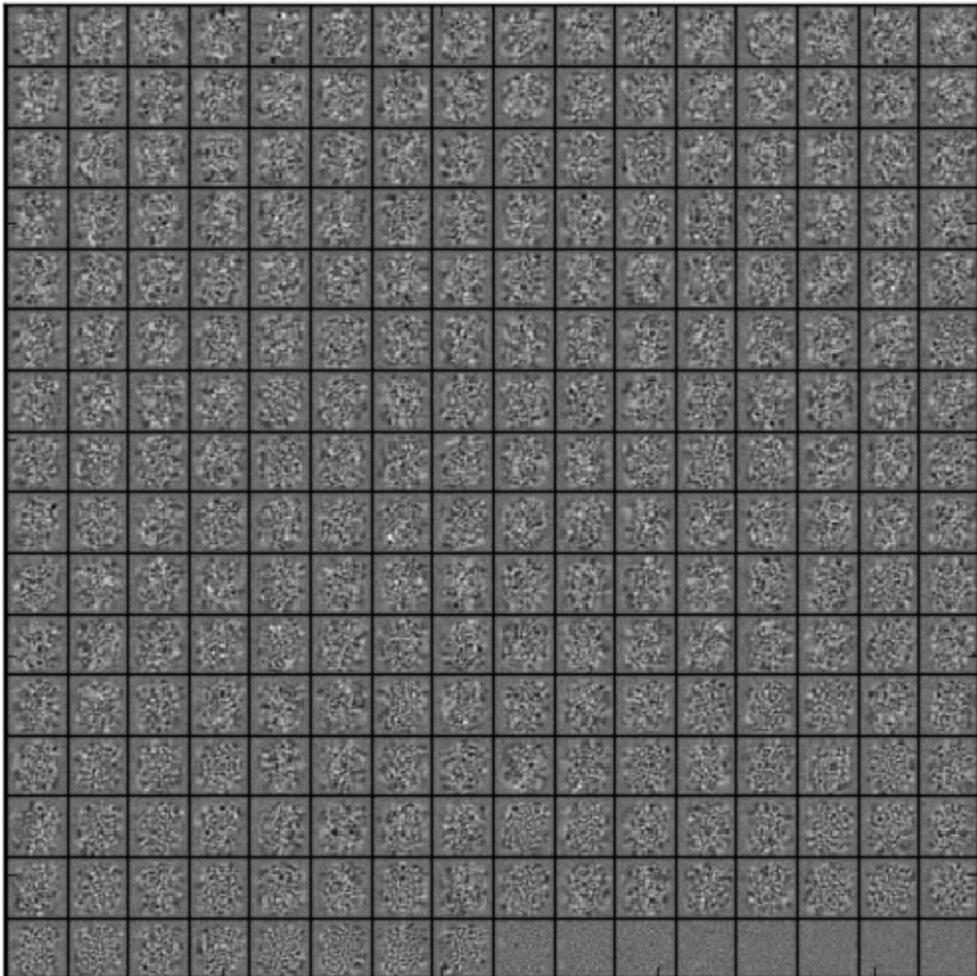


Dropout Results

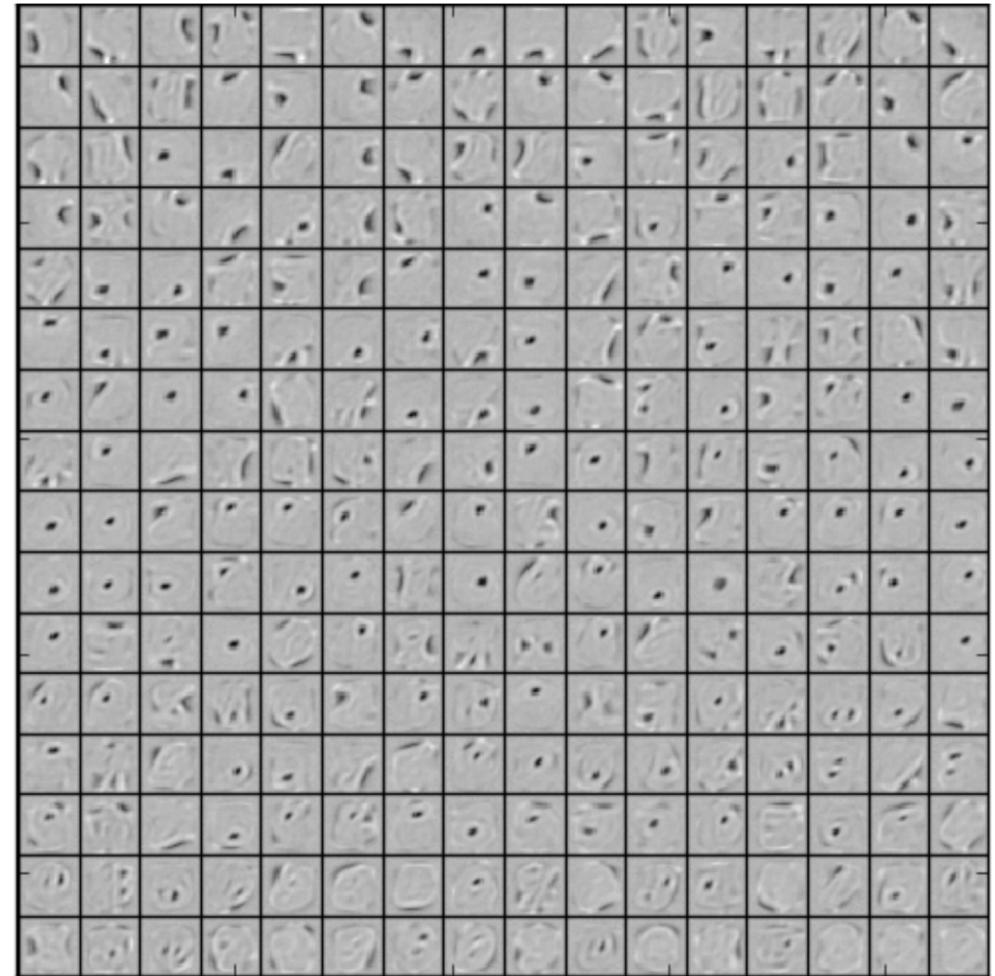


[Srivastava et al JMLR '14]

Dropout: Effect on filters



(a) Without dropout



(b) Dropout with $p = 0.5$.

Batch normalization

- We have discussed data normalization (mean 0, variance 1) as a common approach for stabilizing training
- In deep learning, inputs are shifted and scaled through each layer
- **Batch Normalization** is a widely used technique that normalizes inputs to each layer according to mini-batch statistics
 - “Reduces internal covariate shift”
 - Enables larger learning rates
 - Helps with regularization

Batch normalization [Ioffe & Szegedy 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\underline{\mu}_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\underline{\sigma}_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\underline{\hat{x}_i} \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \underline{\gamma} \hat{x}_i + \underline{\beta} \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

$$\Psi(wx) \rightarrow \Psi(w \text{BN}_{\gamma, \beta}(x))$$

↓
Param: w ↪ Param $\Psi(w, \gamma, \beta)$

Batch normalization [Ioffe & Szegedy 2015]

