# Machine Learning from Data

Lecture 22: Spring 2021

# Today's Lecture

- Neural Networks and Overfitting
    - Approximation Vs Generalization
    - Regularization and Early Stopping
    - Minimizing in-sample error more efficiently

# Neural Networks and Fitting the Data

**Forward Propagation:**

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{\mathbf{W}^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{\mathbf{W}^{(2)}} \mathbf{s}^{(2)} \cdots \xrightarrow{\mathbf{W}^{(L)}} \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x})$$

$$\mathbf{s}^{(\ell)} = (\mathbf{W}^{(\ell)})^{\mathsf{T}} \mathbf{x}^{(\ell-1)} \qquad \mathbf{x}^{(\ell)} = \begin{bmatrix} 1 \\ \theta(\mathbf{s}^{(\ell)}) \end{bmatrix}$$

(Compute $h$ and $E_{\text{in}}$)

Choose $\mathbf{W} = \{\mathbf{W}^{(1)}, \mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}\}$ to minimize $E_{\text{in}}$

**Gradient descent:**

$$\mathbf{W}(t+1) \leftarrow \mathbf{W}(t) - \eta \nabla E_{\text{in}}(\mathbf{W}(t))$$

Compute gradient $\longrightarrow$ need $\dfrac{\partial \mathbf{e}}{\partial \mathbf{W}^{(\ell)}} \longrightarrow$ need $\boldsymbol{\delta}^{(\ell)} = \dfrac{\partial \mathbf{e}}{\partial \mathbf{s}^{(\ell)}}$

$$\frac{\partial \mathbf{e}}{\partial \mathbf{W}^{(\ell)}} = \mathbf{x}^{(\ell-1)} (\boldsymbol{\delta}^{(\ell)})^{\mathsf{T}}$$

**Backpropagation:**

$$\boldsymbol{\delta}^{(1)} \longleftarrow \boldsymbol{\delta}^{(2)} \cdots \longleftarrow \boldsymbol{\delta}^{(L-1)} \longleftarrow \boldsymbol{\delta}^{(L)}$$

$$\boldsymbol{\delta}^{(\ell)} = \theta'(\mathbf{s}^{(\ell)}) \otimes \left[ \mathbf{W}^{(\ell+1)} \boldsymbol{\delta}^{(\ell+1)} \right]_{1}^{d^{(\ell)}}$$



Gradient Descent

SGD

$\log_{10}(\text{error})$

$\log_{10}(\text{iteration})$



Symmetry

Average Intensity

# 2-Layer Neural Network



Perceptrons

$$w_1 \, \theta \left( V_1^T x \right)$$
$$w_2 \, \theta \left( V_2^T x \right)$$
$$\vdots$$
$$w_m \, \theta \left( V_m^T x \right)$$

$d$

m-nodes

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1}^{m} w_j \theta \left( \mathbf{v}_j^{\mathrm{T}} \mathbf{x} \right) \right)$$

# The Neural Network has a Tunable Transform

**Linear Models.**

$$\Phi(x) = [1, \phi_1(x), \phi_2(x) \ldots \phi_m(x)]$$

$k = m$

### Neural Network

$$h(\mathbf{x}) = \theta\left(w_0 + \sum_{j=1}^{m} w_j \theta(\mathbf{v}_j^{\mathsf{T}}\mathbf{x})\right)$$

feature transform fits the data.

### Nonlinear Transform

$$h(\mathbf{x}) = \theta\left(w_0 + \sum_{j=1}^{\tilde{d}} w_j \Phi_j(\mathbf{x})\right)$$

$q$ $m$

$Z$ $E_{in} \times E_{out}$

### k-RBF-Network

$$h(\mathbf{x}) = \theta\left(w_0 + \sum_{j=1}^{k} w_j \phi(\|\mathbf{x} - \boldsymbol{\mu}_j\|)\right)$$
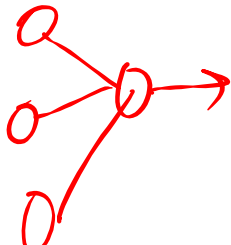
kernel

non-linear

Similarity $\boxed{\mu_j's}$

$$E_{\text{in}} = O\left(\frac{1}{m}\right)$$

approximation

Regression

$$E_{in} = O\left(\frac{1}{\sqrt{m}}\right)$$

$\theta = \text{sign}$

$h(x) = \theta\left(w_0 + \sum_{j=1}^{m} w_j \, \theta(v_j^T x)\right)$

$O(md)$

# Generalization

MLP:

$$d_{\text{VC}} = O(md \log(md))$$

$$m = \sqrt{N}$$

(convergence to optimal for MLP, just like $k$-NN)

*semi-parametric* because you still have to learn parameters.

tanh :

$$d_{\text{VC}} = O(md(m + d))$$

$$E_{in} \in O\left(\frac{1}{m}\right)$$

$$E_{in} \in O\left(\frac{1}{\sqrt{m}}\right)$$

$$E_{out} \leq E_{in} + O\left(\sqrt{\frac{d_{vc}}{N} \log N}\right)^{N^{1/4}}, \quad m = \sqrt{N}$$

$$E_{in} \simeq \frac{1}{N^{1/4}} \qquad E_{out}^{*} \qquad m = \sqrt{N}$$

$$E_{out} \simeq 0 \qquad \boxed{m}$$

Semi-parametric method.    Validation

$$x \longrightarrow \left( h_1(x) \quad h_2(x) \dots h_m(x) \right)$$

$$h(x) = h_c \left( h_1(x) \quad h_2(x) \right.$$
$$\left. \dots h_m(x) \right)$$

$$\longrightarrow h(x) \quad \# \text{ of dichotomies}$$

$$x_1 \longrightarrow \begin{pmatrix} h_1(x_1) & h_2(x_1) \cdots h_m(x_1) \end{pmatrix} \longrightarrow z_1$$

$$x_2 \longrightarrow \begin{pmatrix} h_1(x_2) & h_2(x_2) \cdots h_m(x_2) \end{pmatrix} \longrightarrow z_2$$

$$\vdots$$

$$x_N \longrightarrow \begin{pmatrix} h_1(x_N) & h_2(x_N) \cdots h_m(x_N) \end{pmatrix} \longrightarrow z_m$$

$$N^{d_1} \quad N^{d_2} \quad - - - \quad N^{d_m}$$

$$\# \text{ of configurations of } z\text{'s} \leq N^{d_1 + d_2 \cdots + d_m}$$

$$\# \text{ of dichotomies} \leq N^{\underbrace{d_c + \sum_{i=1}^{m} d_i}_{D}}$$

$$\text{Lemma:} \quad N^D \leq 2^N \quad , \quad N \in \Omega(D \log D)$$

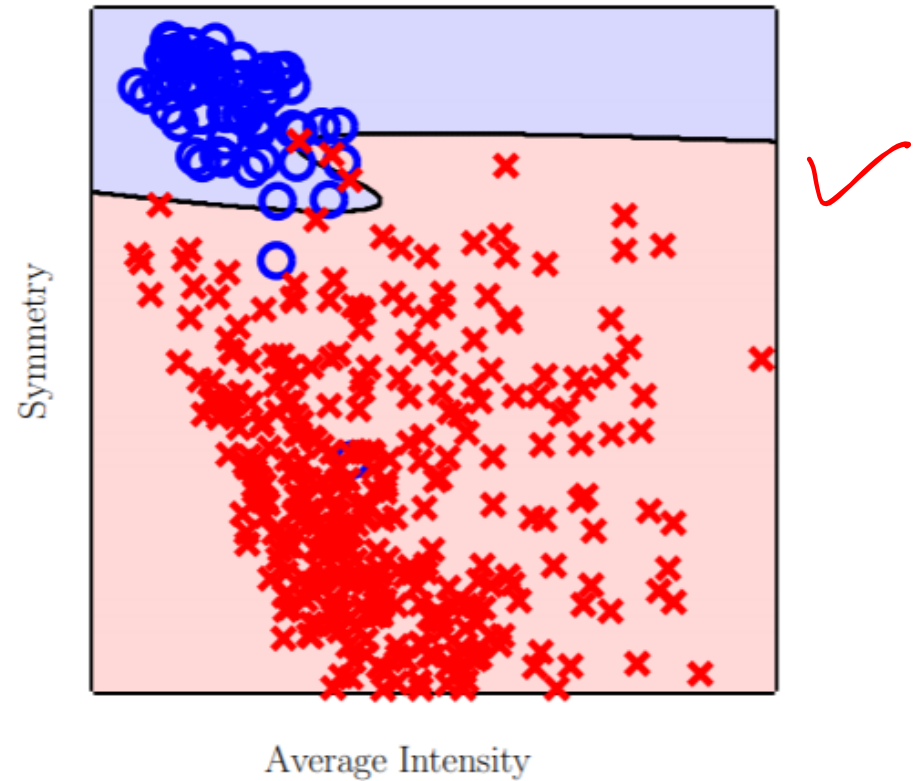$$N > 2D \log_2 D \quad \text{then} \quad m(N) < 2^N \quad \& \quad d_{VC} \leq 2D \log D$$

$$O\left(\sum d_i + d_c\right) \log\left(\sum d_i + d_c\right)$$

# Weight Decay with Digits Data



No Weight Decay

Weight Decay, $\lambda = 0.01$

Symmetry

Average Intensity

Symmetry

Average Intensity

$$E_{aug} = E_{in} + \frac{\lambda}{N} \omega^T \omega$$

$$W = \{ W^{(1)}, W^{(2)} \cdots W^{(l)} \}$$

$$E_{aug} = E_{in} + \frac{\lambda}{N} \sum_{l,i,j} \left( \omega_{ij}^{(l)} \right)^2$$

$$\frac{\partial E_{aug}}{\partial W_{ij}^{(l)}} = \frac{\partial E_{in}}{\partial \omega_{ij}^{(l)}} + \frac{\lambda}{N} 2 \cdot \omega_{ij}^{(l)}$$

$$\rightarrow \frac{\partial E_{aug}}{\partial W^{(l)}} = \frac{\partial E_{in}}{\partial W^{(l)}} + \frac{2\lambda}{N} W^{(l)}$$

Back propagation

$$W^{(l)} \leftarrow W^{(l)} - \eta \frac{\partial E_{in}}{\partial W^{(l)}} - \frac{2\lambda}{N} W^{(l)}$$

Hypothesis
set

$$H_1 \subseteq H_2 \subseteq H_3 \cdots \subseteq H_m$$

$$d_1 \leq d_2 \leq d_3 \cdots \leq d_m.$$

# Early Stopping

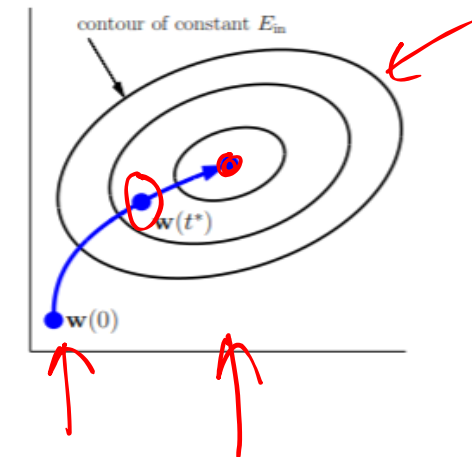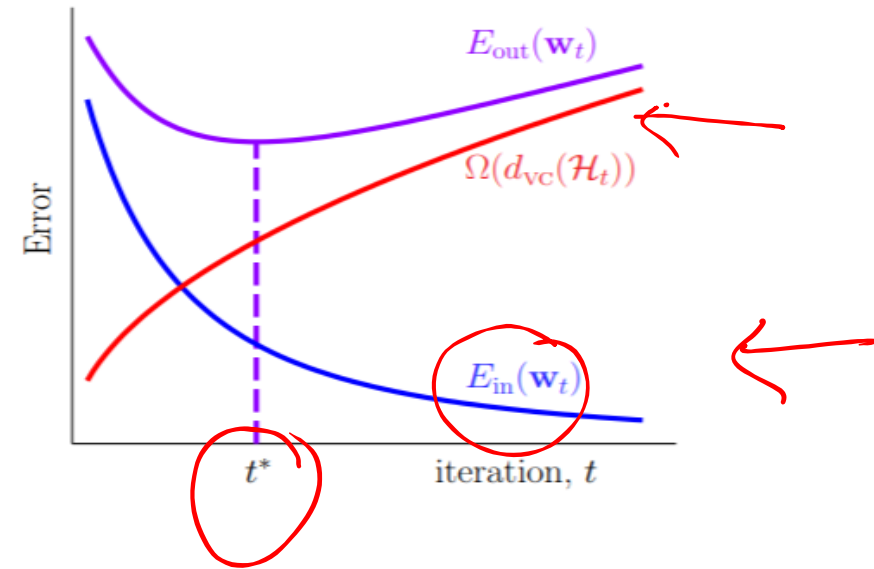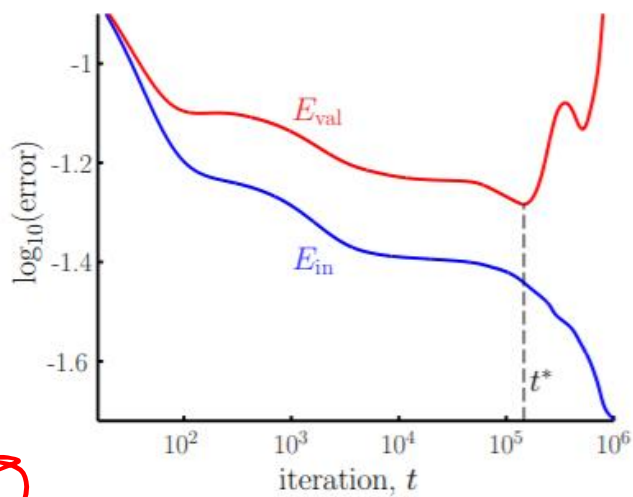## Gradient Descent

$$\mathbf{w}_1 = w_0 - \eta \frac{\mathbf{g}_0}{\|\mathbf{g}_0\|}$$

$$\mathcal{H}_1 = \{\mathbf{w} : \ \|\mathbf{w} - \mathbf{w}_0\| \leq \eta\}$$

$$\mathcal{H}_2 = \mathcal{H}_1 \cup \{\mathbf{w} : \ \|\mathbf{w} - \mathbf{w}_1\| \leq \eta\}$$

$$\mathcal{H}_3 = \mathcal{H}_2 \cup \{\mathbf{w} : \ \|\mathbf{w} - \mathbf{w}_2\| \leq \eta\}$$

Each iteration explores a larger $\mathcal{H}$

$$\mathcal{H}_1 \subset \mathcal{H}_2 \subset \mathcal{H}_3 \subset \mathcal{H}_4 \subset \cdots$$

$E_{\text{out}}(\mathbf{w}_t)$

$\Omega(d_{\text{vc}}(\mathcal{H}_t))$

$E_{\text{in}}(\mathbf{w}_t)$

Error

$t^*$

iteration, $t$

contour of constant $E_{\text{in}}$

$\mathbf{w}(t^*)$

$\mathbf{w}(0)$

# Early Stopping on Digits Data



Use a validation set to determine $t^*$

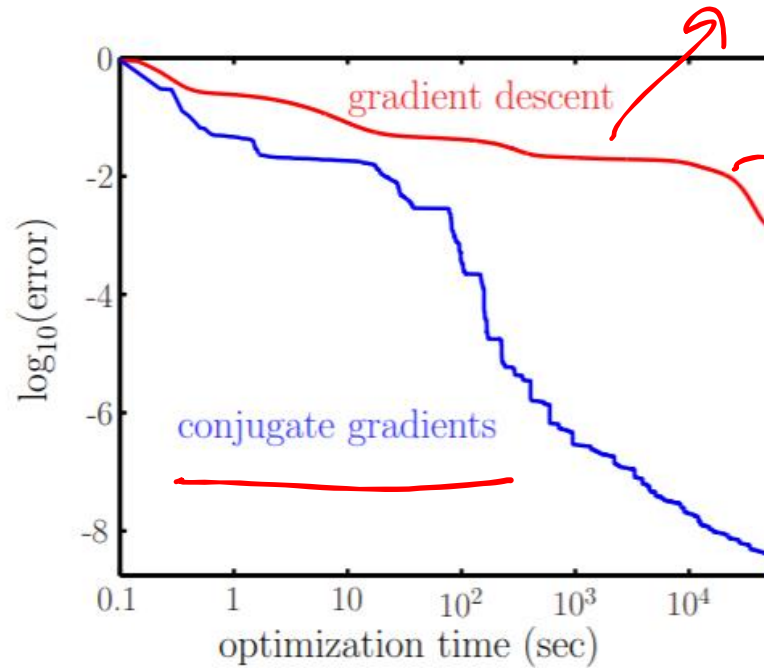Output $\mathbf{w}^*$, do not retrain with all the data till $t^*$.

# Minimizing $E_{in}$

1. Use regression for classification

2. Use better algorithms than gradient descent



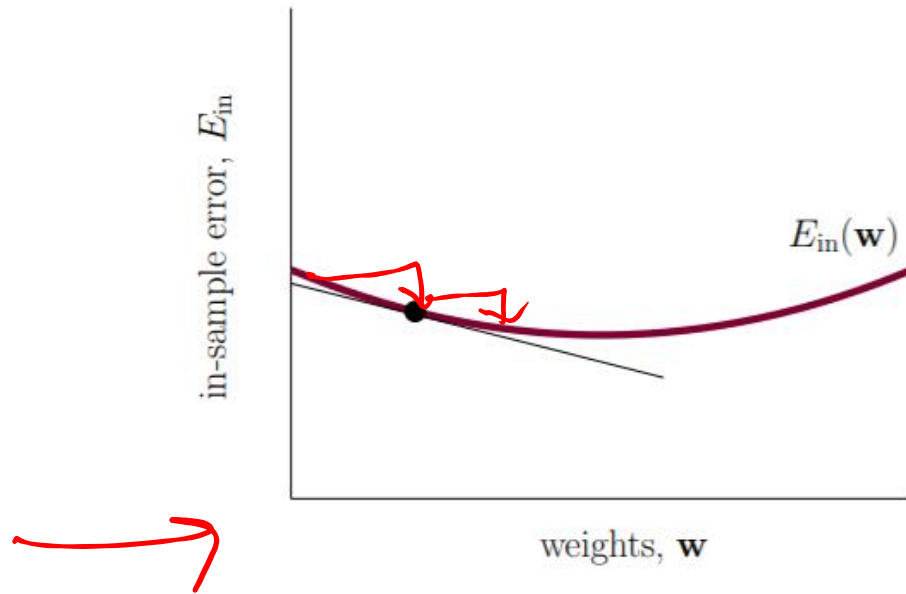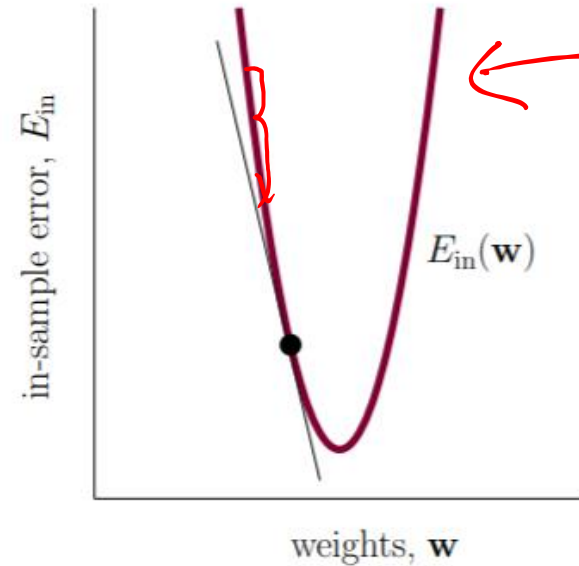$$s^L \longrightarrow \bigcirc \longrightarrow h(x) = \theta(s^L)$$

sign or tanh

fast

i) Direction (-ve of the grad.)
ii) Step size ($\eta$)

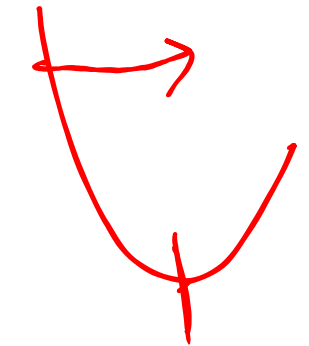Determine the gradient **g**
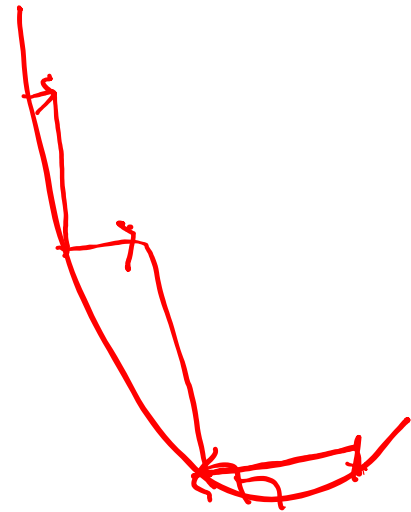


Shallow: use large $\eta$.

Deep: use small $\eta$.

1) Start, take a step, check if $E_{in}$ improves.

$$\eta \longrightarrow \eta \times \text{some factor} (\beta = 1.03)$$

2) Take a bigger $\dfrac{step}{\cdot}$ & improve $E_{in}$.

3) If you worsen $E_{in}$ then
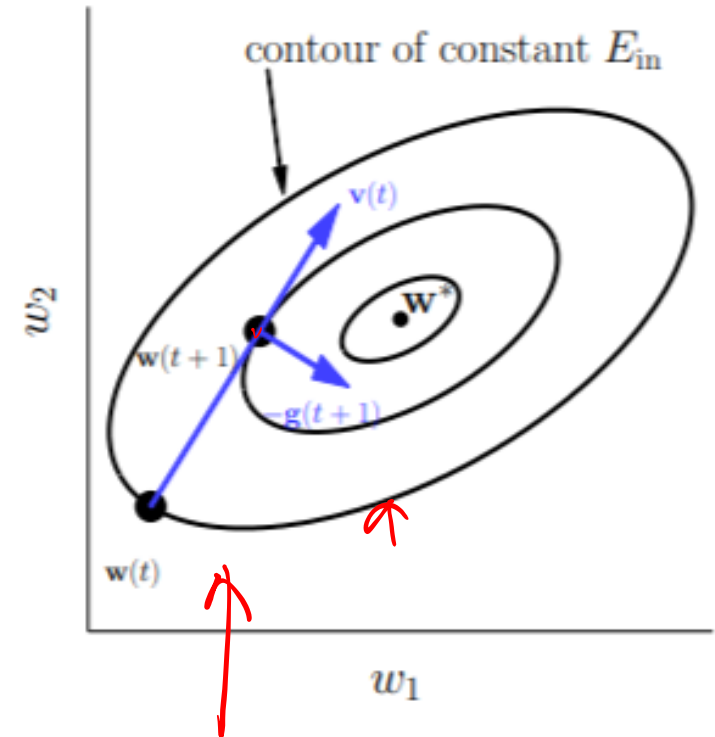
$$\eta \longrightarrow \eta \times \alpha (<1)$$

You may backtrack.

Variable $\eta \longrightarrow$ GD
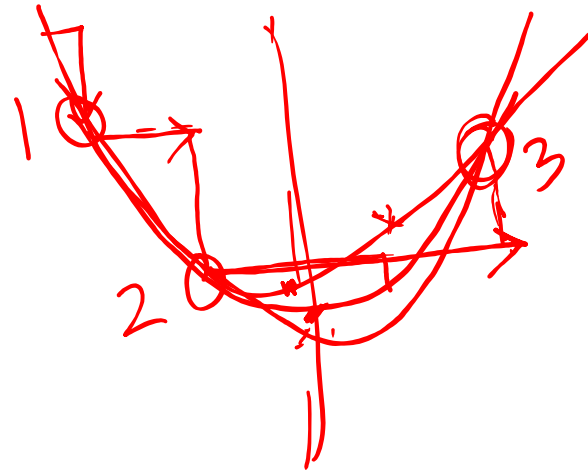
# Steepest Descent - Line Search

1: Initialize $w(0)$ and set $t = 0$;

2: **while** stopping criterion has not been met **do**

3:     Let $\mathbf{g}(t) = \nabla E_{\text{in}}(\mathbf{w}(t))$, and set $\mathbf{v}(t) = -\mathbf{g}(t)$.

4:     Let $\eta^* = \text{argmin}_\eta E_{\text{in}}(\mathbf{w}(t) + \eta\mathbf{v}(t))$.

5:     $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta^*\mathbf{v}(t)$.

6:     Iterate to the next step, $t \leftarrow t + 1$.

7: **end while**



How to accomplish the line search (step 4)?

Simple bisection (binary search) suffices in practice

— U shaped

Binary search

In practice 5-10 ($\varepsilon$)

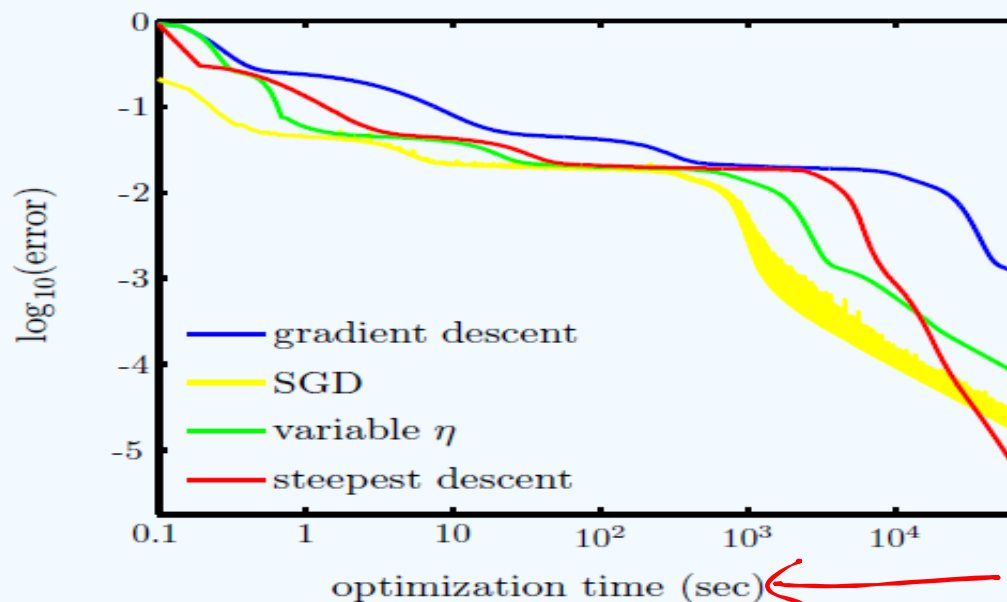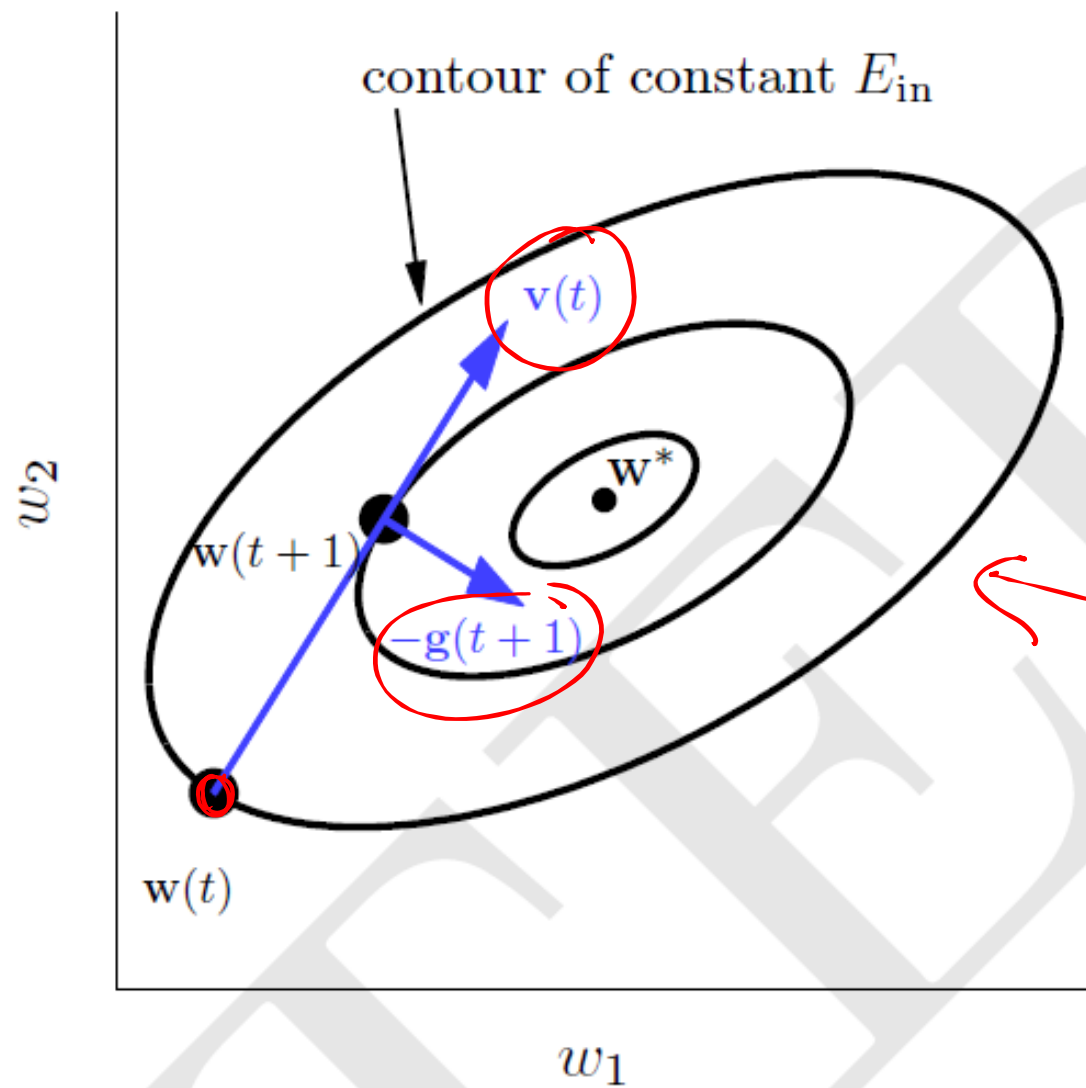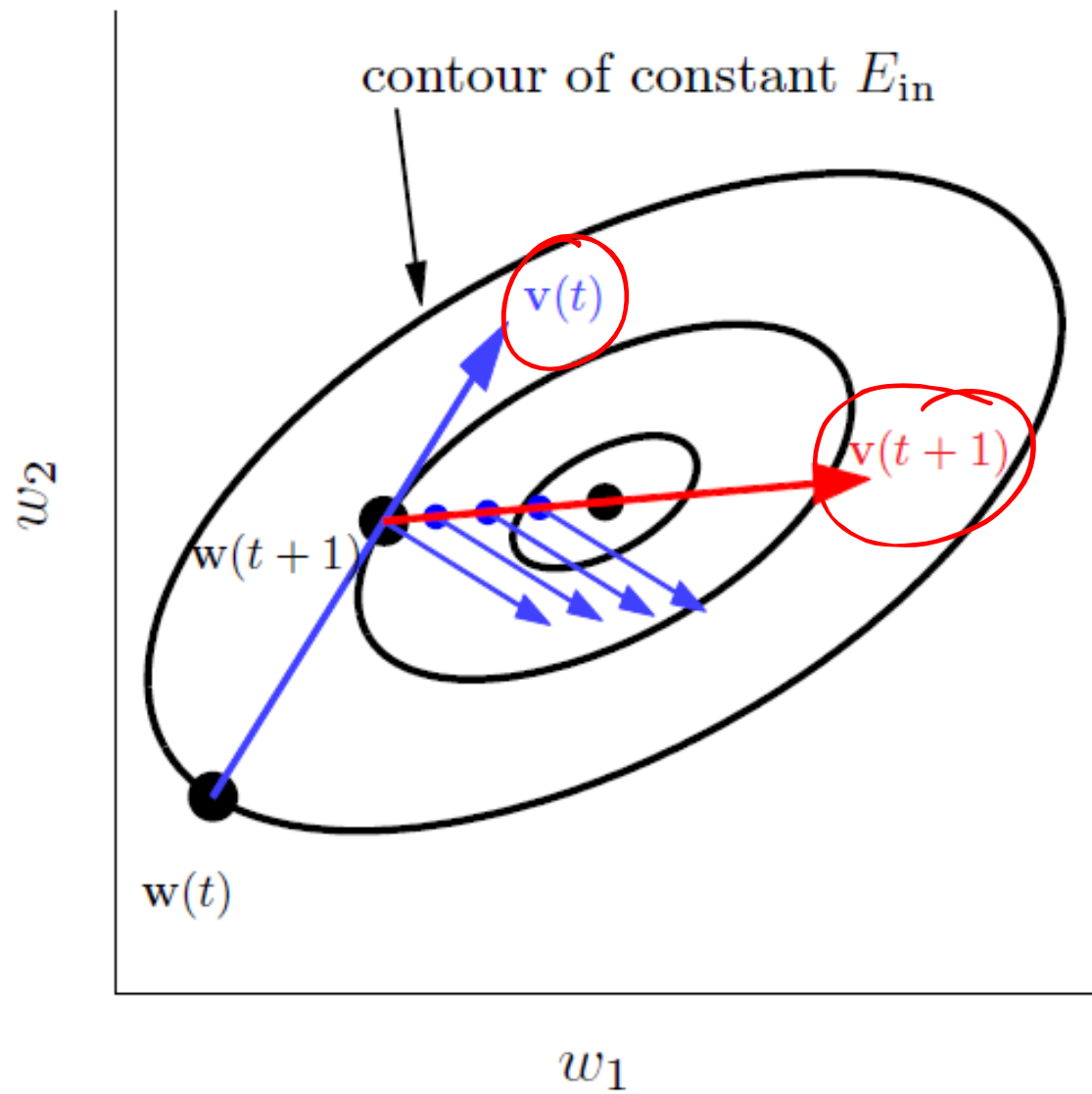| | Optimization Time | | |
| Method | 10 sec | 1,000 sec | 50,000 sec |
|---|---|---|---|
| Gradient Descent | 0.079 | 0.0206 | 0.00144 |
| **Stochastic Gradient Descent** | **0.0213** | **0.00278** | 0.000022 |
| Variable Learning Rate | 0.039 | 0.014 | 0.00010 |
| Steepest Descent | 0.043 | 0.0189 | **0.000012** |



Figure 7.4: Gradient descent, variable learning rate and steepest descent using digits data and a 5 hidden unit 2-layer neural network with linear output. For variable learning rate, $\alpha = 1.1$ and $\beta = 0.8$.

contour of constant $E_{in}$

$\mathbf{v}(t)$

$\mathbf{w}^*$

$\mathbf{w}(t+1)$

$-\mathbf{g}(t+1)$

$\mathbf{w}(t)$

$w_2$

$w_1$

new grad ⊥ prev. line search dir

contour of constant $E_{\text{in}}$

$\mathbf{v}(t)$
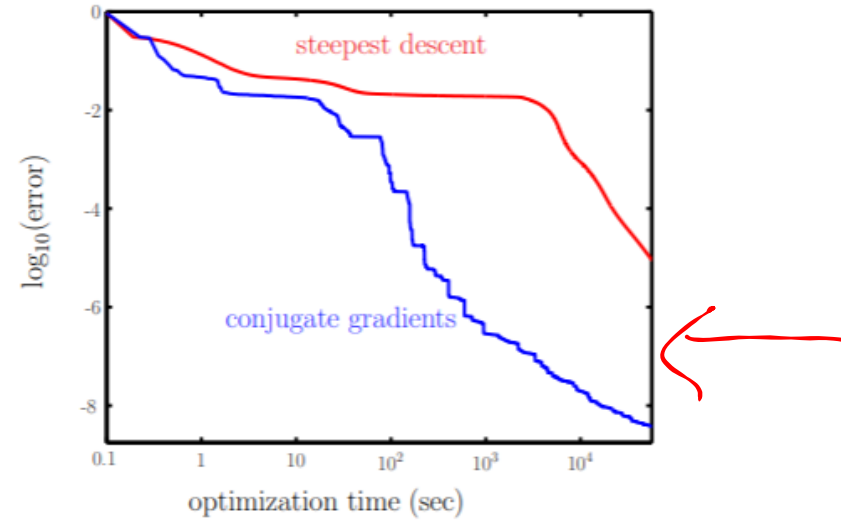
$\mathbf{v}(t+1)$

$\mathbf{w}(t+1)$

$\mathbf{w}(t)$

$w_2$

$w_1$

$\alpha$-steps.
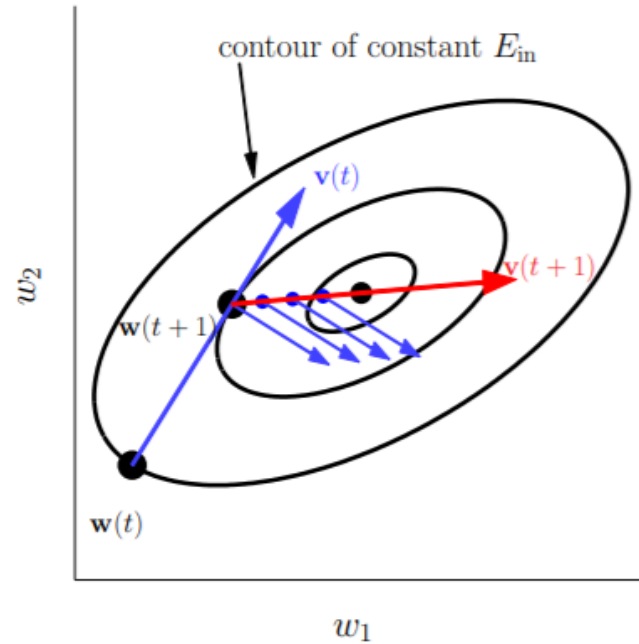
# Conjugate Gradients

1. Line search just like steepest descent.

2. Choose a better direction than $-\mathbf{g}$



contour of constant $E_{\text{in}}$

| | Optimization Time | | |
|---|---|---|---|
| Method | 10 sec | 1,000 sec | 50,000 sec |
| Stochastic Gradient Descent | 0.0203 | 0.000447 | $1.6310 \times 10^{-5}$ |
| Steepest Descent | 0.0497 | 0.0194 | 0.000140 |
| **Conjugate Gradients** | **0.0200** | $\mathbf{1.13 \times 10^{-6}}$ | $\mathbf{2.73 \times 10^{-9}}$ |

# Thanks!