

Regularization in Neural Networks

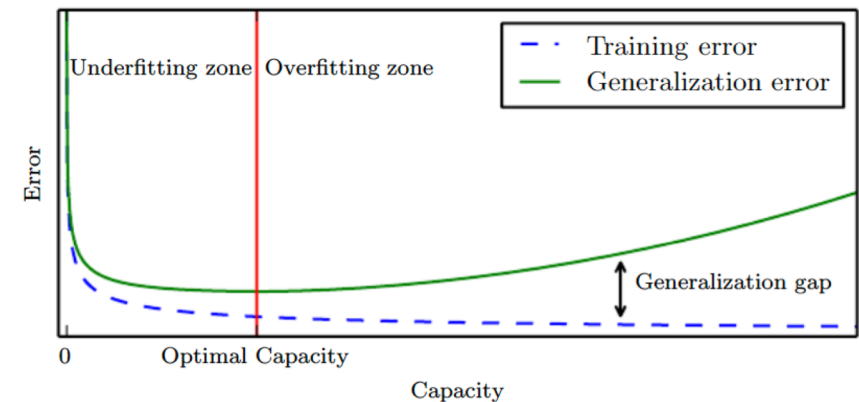
Sargur Srihari
srihari@buffalo.edu

Topics in Neural Net Regularization

- Definition of regularization
- Methods
 1. Limiting capacity: no of hidden units
 2. Norm Penalties: L^2 - and L^1 - regularization
 3. Early stopping
 - Invariant methods
 4. Data Set Augmentation
 5. Tangent propagation

Definition of Regularization

- Generalization error
 - Performance on inputs not previously seen
 - Also called as *Test error*
- Regularization is:
 - “any modification to a learning algorithm to reduce its *generalization error* but not its *training error*”
 - Reduce generalization error even at the expense of increasing training error
 - E.g., Limiting model capacity is a regularization method



Goals of Regularization

- Some goals of regularization
 1. Encode prior knowledge
 2. Express preference for simpler model
 3. Need to make underdetermined problem determined

Need for Regularization

- Generalization
 - Prevent over-fitting
- Occam's razor
- Bayesian point of view
 - Regularization corresponds to prior distributions on model parameters

Best Model: one properly regularized

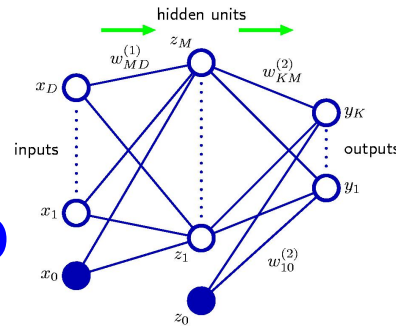
- Best fitting model obtained not by finding the right number of parameters
- Instead, best fitting model is a large model that has been regularized appropriately
- We review several strategies for how to create such a large, deep regularized model

Limiting no. of hidden units

- No. of input/output units determined by dimensions
- Number of hidden units M is a free parameter
 - Adjusted to get best predictive performance
- Possible approach is to get maximum likelihood estimate of M for balance between under-fitting and over-fitting

Effect of No. of Hidden Units

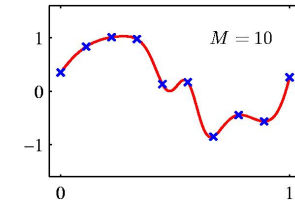
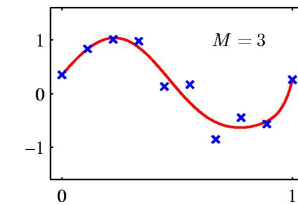
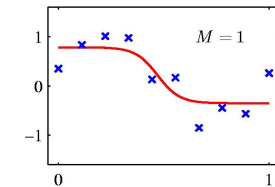
Sinusoidal Regression Prob



Two layer network trained on 10 data points

$M = 1, 3$ and 10 hidden units

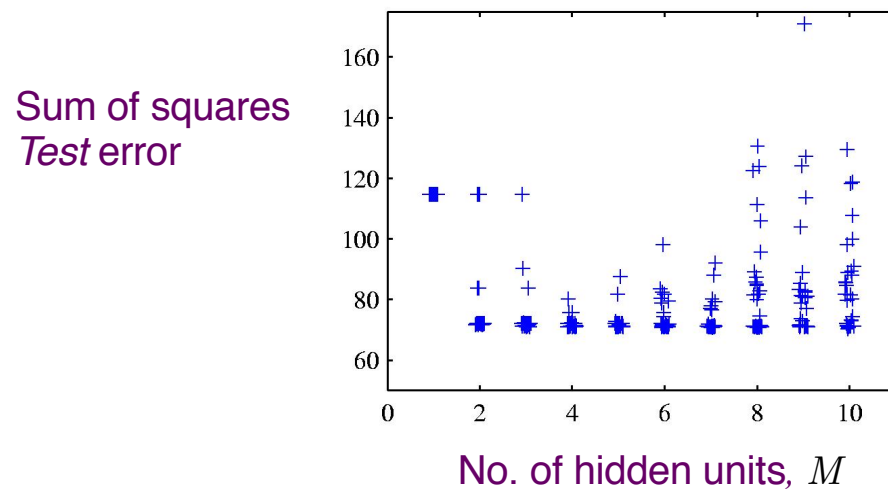
Minimizing sum-of-squared error function
Using conjugate gradient descent



Generalization error is not a simple function of M
due to presence of local minima in error function

Validation Set to fix no. of hidden units

- Plot a graph choosing random starts and different numbers of hidden units M and choose the specific solution having smallest generalization error



- 30 random starts for each M
30 points in each column of graph
- Overall best *validation set* performance happened at $M=8$

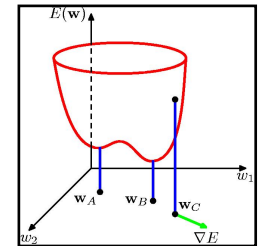
- There are other ways to control the complexity of a neural network in order to avoid over-fitting
- Alternative approach is to choose a relatively large value of M and then control complexity by adding a regularization term

Parameter Norm Penalty

- Generalization error not a simple function of M
 - Due to presence of local minima
 - Need to control capacity to avoid over-fitting
 - Alternatively choose large M and control complexity by addition of regularization term

- Simplest regularizer is *weight decay*

$$\tilde{E}(w) = E(w) + \frac{\lambda}{2} w^T w$$



- Effective model complexity determined by choice of λ
- Regularizer is equivalent to a Gaussian prior over w
- In a 2-layer neural network

$$J_{\text{regularized}} = \underbrace{-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}))}_{\text{cross-entropy cost}} + \underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_l \sum_k \sum_j W_{kj}^{[l]2}}_{\text{L2 regularization cost}}$$

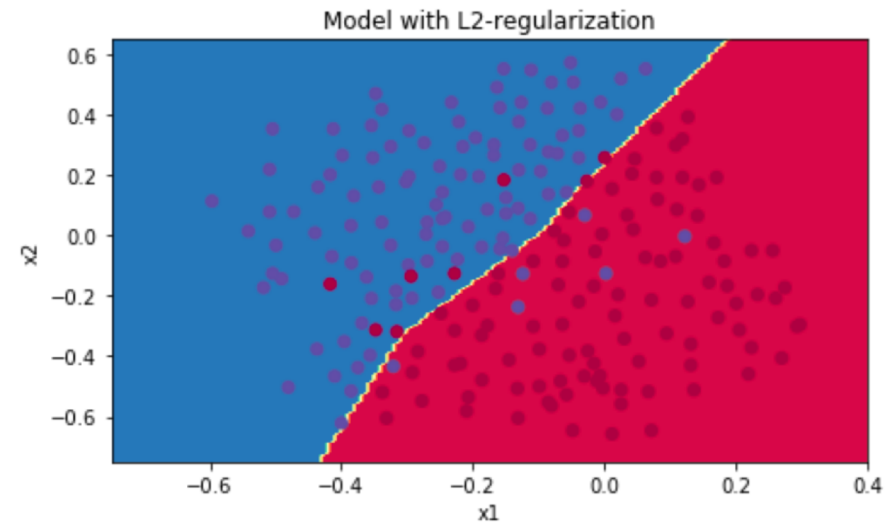
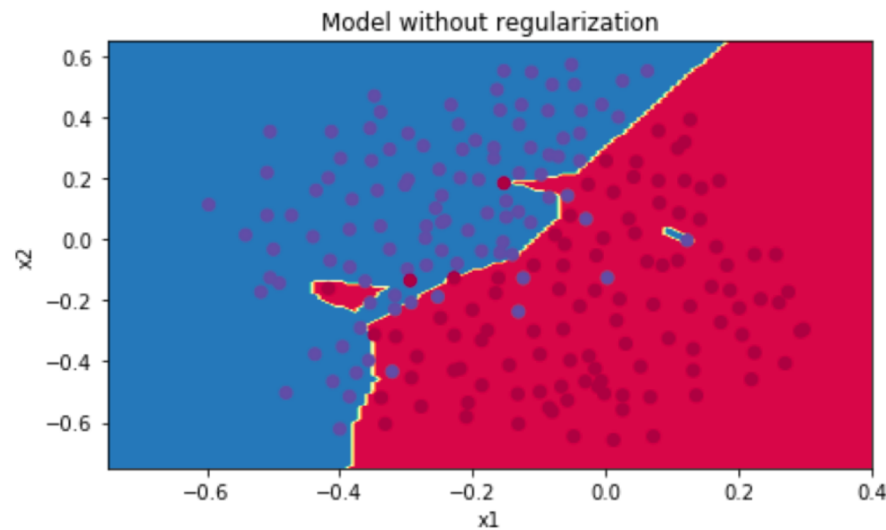
Weight Decay

- The name weight decay is due to the following

$$w_{t+1} = w_t - \alpha \nabla_w J - \lambda w_t$$

- To prevent overfitting, every time we update a weight w with the gradient ∇J in respect to w , we also subtract from it λw .
- This gives the weights a tendency to decay towards zero, hence the name.

Effect of Neural Net Regularization



$$J_{\text{regularized}} = \underbrace{-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}))}_{\text{cross-entropy cost}} + \underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_l \sum_k \sum_j W_{kj}^{[l]2}}_{\text{L2 regularization cost}}$$

https://medium.com/@shay.palachy/understanding-the-scaling-of-l%C2%B2-regularization-in-the-context-of-neural-networks-e3d25f8b50db?fbclid=IwAR17rKfnM2JnkWBOZFvQg6ftl0k52diyEra_66UCwI18Gtf7nFd6l3l4VMY

Invariance to Transformation

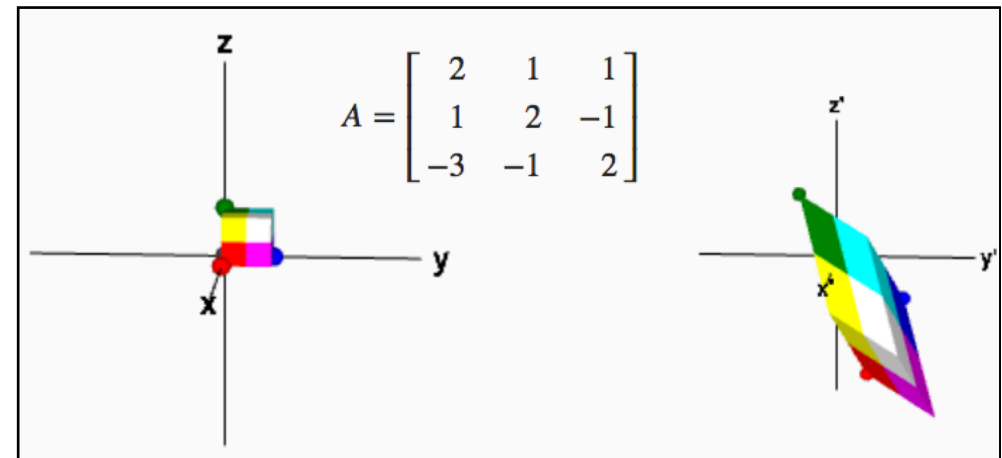
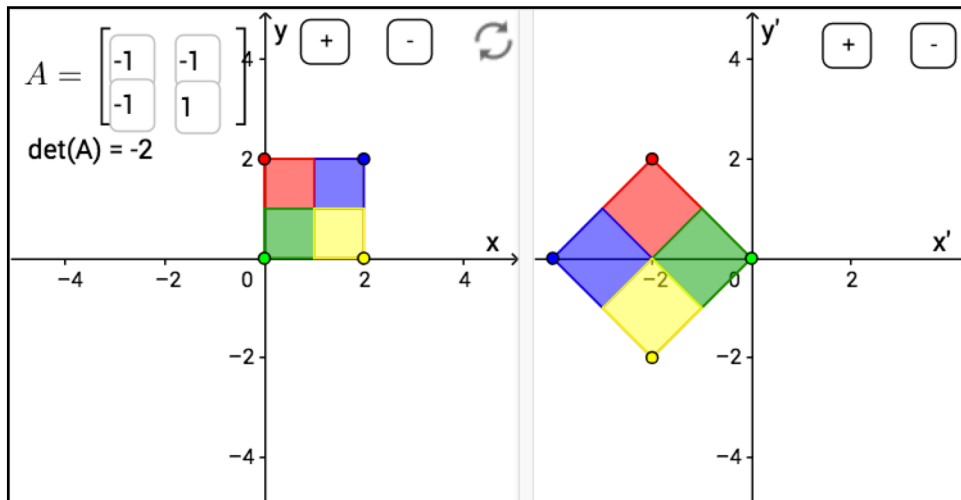
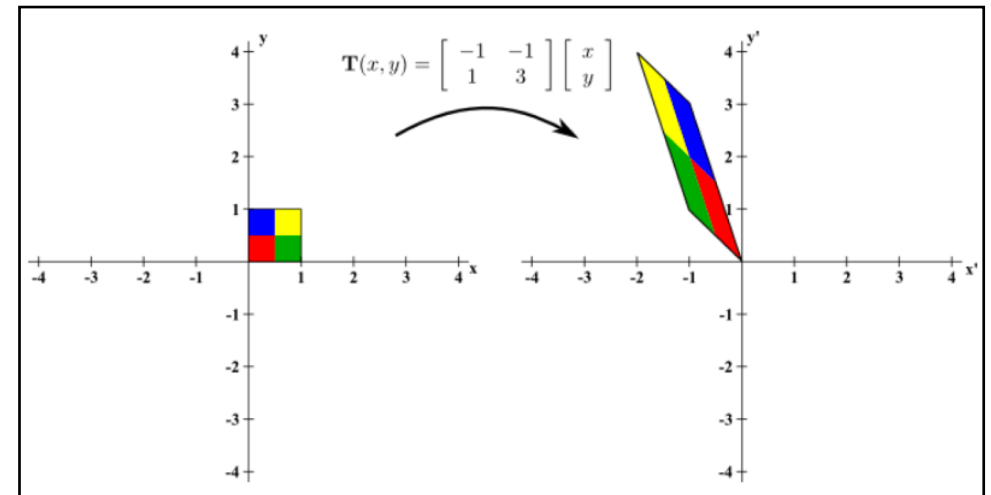
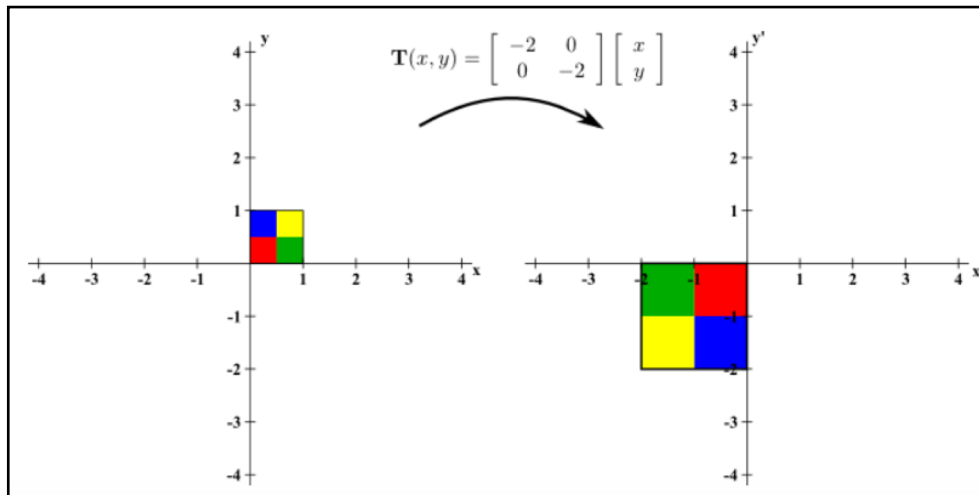
- Simple weight decay has certain shortcomings
invariance to scaling
 - Same building, different views, scales



Linear Transformation

$$\mathbf{T}(x, y) = (ax + by, cx + dy) = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\mathbf{T}(\mathbf{x}) = \mathbf{A}\mathbf{x}$$



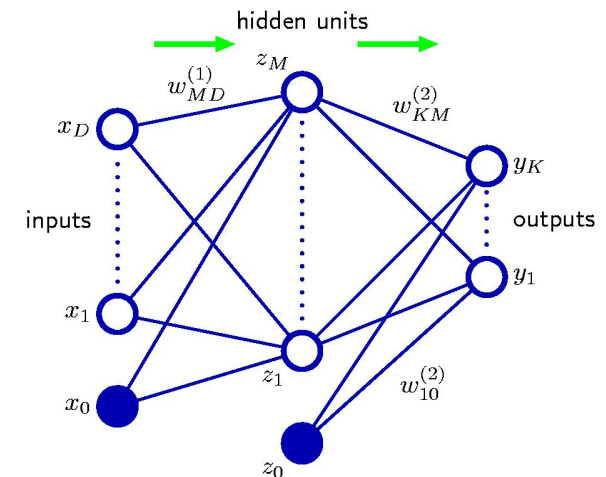
Ex: Two-layer MLP

- Simple weight decay is inconsistent with certain scaling properties of network mappings
- To show this, consider a multi-layer perceptron network with two layers of weights and linear output units
 - Set of inputs $\mathbf{x} = \{x_i\}$ and outputs $\mathbf{y} = \{y_k\}$
 - Activations of hidden units in first layer have form

$$z_j = h\left(\sum_i w_{ji} x_i + w_{j0}\right)$$

- Activations of output units are

$$y_k = \sum_j w_{kj} z_j + w_{k0}$$



Linear Transformations of i/o Variables

- If we perform linear transformation of input data

$$x_i \rightarrow \tilde{x}_i = ax_i + b$$

- Then mapping performed by network is unchanged by making a corresponding linear transformation of
 - weights, biases from the inputs to the hidden units as

$$w_{ji} \rightarrow \tilde{w}_{ji} = \frac{1}{a} w_{ji} \quad \text{and} \quad w_{j0} = w_{j0} - \frac{b}{a} \sum_i w_{ji}$$

- Similar linear transform of outputs of network

$$y_k \rightarrow \tilde{y}_k = cy_k + d$$

- achieved by transforming second layer weights, biases

$$w_{kj} \rightarrow \tilde{w}_{kj} = cw_{kj} \quad \text{and} \quad w_{k0} = cw_{k0} + d$$

Desirable invariance of regularizer

- Suppose we train two different networks
 - *First network*: trained using original data: $\mathbf{x}=\{x_i\}$, $\mathbf{y}=\{y_k\}$
 - *Second network*: input and/or target variables are transformed by one of the linear transformations

$$x_i \rightarrow \tilde{x}_i = ax_i + b$$

$$y_k \rightarrow \tilde{y}_k = cy_k + d$$

- Consistency requires that we should obtain equivalent networks that differ only by linear transformation of the weights

For first layer:

$$w_{ji} \rightarrow \frac{1}{a} w_{ji}$$

and/or

And/or or second layer:

$$w_{kj} \rightarrow cw_{kj}$$

Simple weight decay fails invariance

- Regularizer should have this property
 - Otherwise arbitrarily favor one solution over another
- Simple weight decay $\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$ does not have this property
 - Treats all weights and biases on an equal footing
 - While resulting w_{ji} and w_{kj} should be treated differently
 - Consequently networks will have different weights and violate invariance
 - We therefore look for a regularizer invariant under the linear transformations

Regularizer invariance

- The regularizer should be invariant to re-scaling of weights and shifts of biases
- Such a regularizer is
$$\frac{\lambda_1}{2} \sum_{w \in W_1} w^2 + \frac{\lambda_2}{2} \sum_{w \in W_2} w^2$$
 - where W_1 are weights of first layer and
 - W_2 are the set of weights in the second layer
- This regularizer remains unchanged under the weight transformations provided the parameters are rescaled using $\lambda_1 \rightarrow a^{1/2} \lambda_1$ and $\lambda_2 \rightarrow c^{-1/2} \lambda_2$
- We have seen before that weight decay is equivalent to a Gaussian prior. So what is the equivalent prior to this?

Bayesian linear regression

$$p(\mathbf{t} \mid X, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n \mid \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

$$p(\mathbf{w} \mid \alpha) = N(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I})$$

$$p(\mathbf{w} \mid \mathbf{t}) = \prod_{n=1}^N N(t_n \mid \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) N(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I})$$

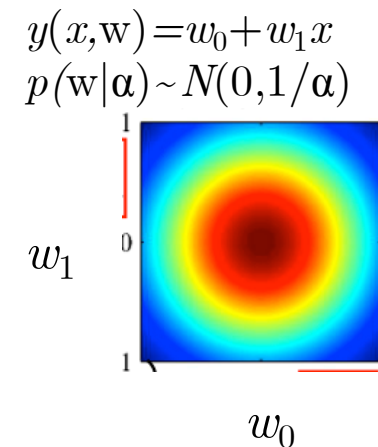
Log of Posterior is

$$\ln p(\mathbf{w} \mid \mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const}$$

Thus Maximization of posterior is equivalent to weight decay

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

with addition of quadratic regularization term $\mathbf{w}^T \mathbf{w}$ with $\lambda = \alpha / \beta$



Equivalent prior for neural network

- Regularizer invariant to rescaled weights/ biases

$$\frac{\lambda_1}{2} \sum_{w \in W_1} w^2 + \frac{\lambda_2}{2} \sum_{w \in W_2} w^2$$

- where W_1 : weights of first layer, W_2 : of second layer
- It corresponds to a prior of the form

$$p(\mathbf{w} \mid \alpha_1, \alpha_2) \propto \exp \left(-\frac{\alpha_1}{2} \sum_{w \in W_1} w^2 - \frac{\alpha_2}{2} \sum_{w \in W_2} w^2 \right) \quad \alpha_1 \text{ and } \alpha_2 \text{ are hyper-parameters}$$

- An *improper prior* which cannot be normalized
 - Leads to difficulties in selecting regularization coefficients and in model comparison within Bayesian framework
 - Instead include separate priors for biases with their own hyper-parameters

Separate priors and biases

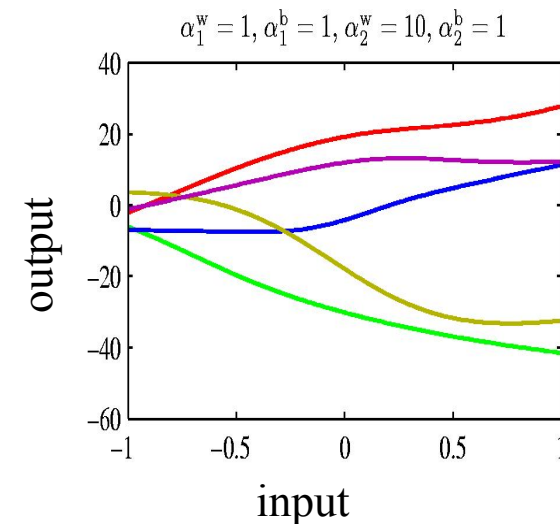
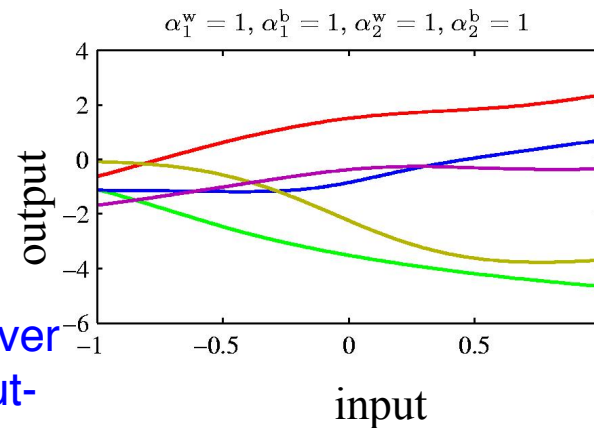
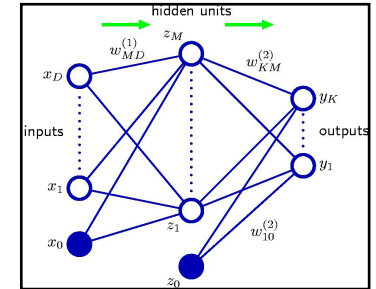
- We can illustrate effect of the resulting four hyperparameters
 - α_1^b , precision of Gaussian distribution of *first layer bias*
 - α_1^w , precision of Gaussian of *first layer weights*
 - α_2^b , precision of Gaussian distribution of *second layer bias*
 - α_2^w , precision of Gaussian distribution of *second layer weights*

Effect of hyperparameters on i/o

Network with single input
(x value range: -1 to +1),
single linear output
(y value range: -60 to +40)

12 hidden units with tanh
activation functions

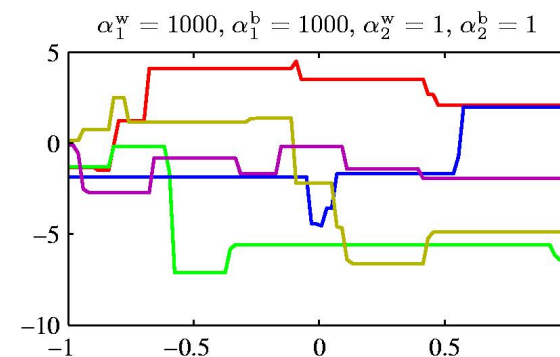
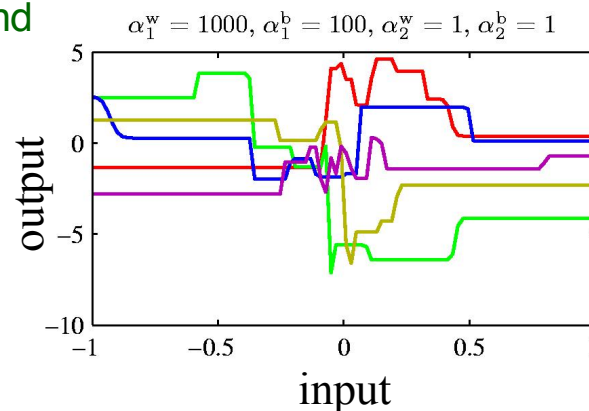
Priors are governed by four hyper-parameters
 α_1^b , precision of Gaussian distribution of *first layer bias*
 α_1^w , precision of Gaussian of *first layer weights*
 α_2^b, \dots of *second layer bias*
 α_2^w, \dots of *second layer weights*



Observe:
vertical
scale
controlled
by α_2^w

Draw five samples from prior over
hyperparameters, and plot input-
output (network functions)

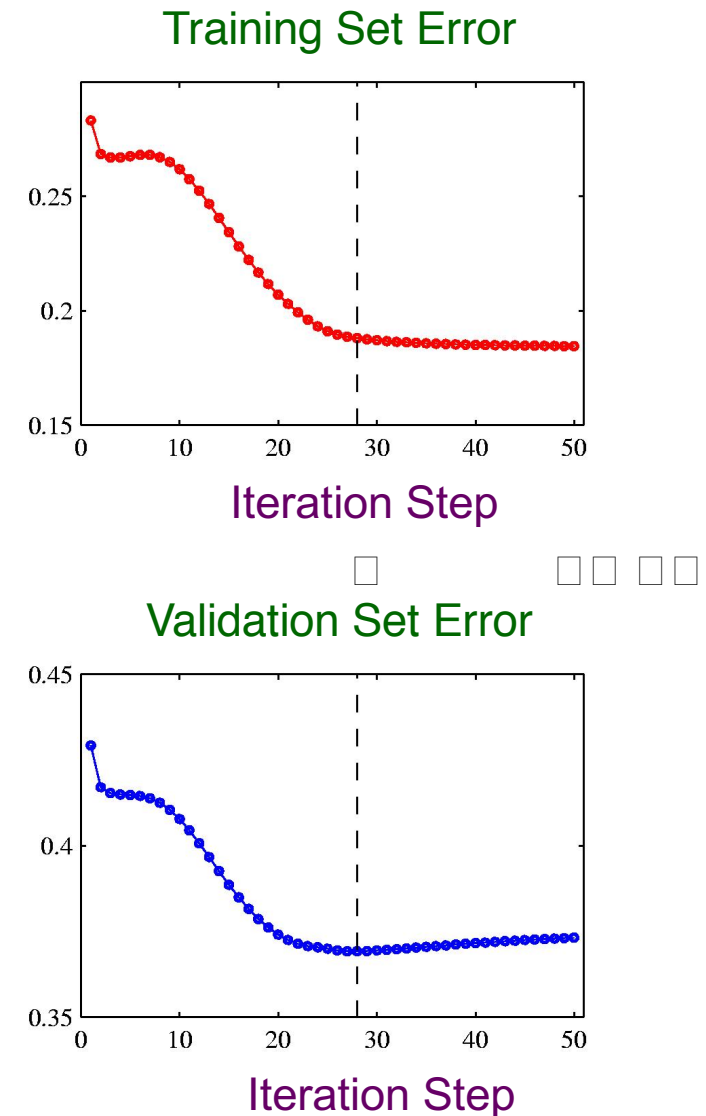
Five samples correspond to five colors
For each setting function is learnt and
plotted



Observe:
horizontal
scale
controlled
by α_1^w

3. Early Stopping

- Alternative to regularization
 - In controlling complexity
- Error measured with an independent validation set
 - shows initial decrease in error and then an increase
- Training stopped at point of smallest error with validation data
 - Effectively limits network complexity

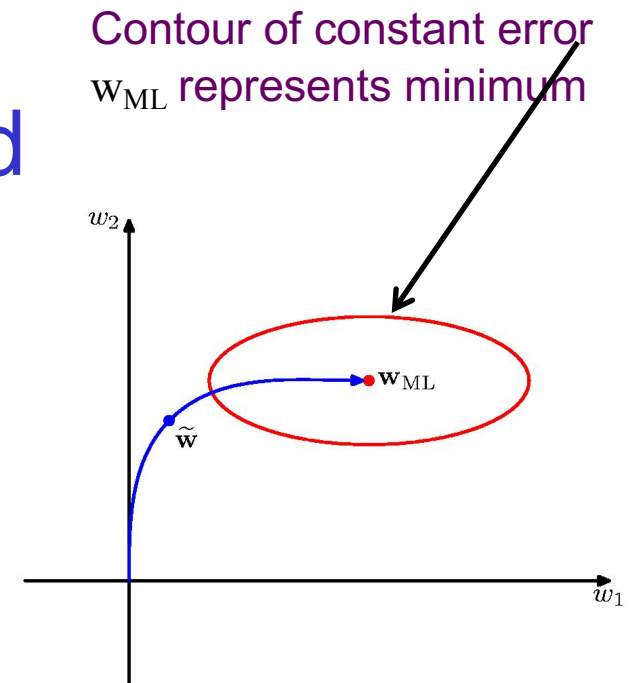


Interpreting the effect of Early Stopping

- Consider quadratic error function
 - Axes in weight space are parallel to eigen vectors of Hessian
- In absence of weight decay, weight vector starts at origin and proceeds to \mathbf{w}_{ML}
- Stopping at $\tilde{\mathbf{w}}$ similar to weight decay

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Effective no of parameters in network grows during training



Invariances

- In classification there is a need for:
 - Predictions invariant to transformations of input
- Example:
 - handwritten digit should be assigned same classification irrespective of position in image (translation) and size (scale)
- Such transformations produce significant changes in raw data, yet need to produce same output from classifier
 - Examples: pixel intensities,
in speech recognition, nonlinear time warping along time axis

Simple Approach for Invariance

- Large sample set where all transformations are present
 - E.g., for translation invariance, examples of objects in many different positions
- Impractical
 - Number of different samples grows exponentially with number of transformations
- Seek alternative approaches for adaptive model to exhibit required invariances

Approaches to Invariance

1. Training set augmented

- Transform training patterns for desired invariances
E.g., shift each image into different positions

2. Regularization term to error function

- Penalize changes in output when input transformed
 - Leads to tangent propagation.

3. Invariance built into pre-processing

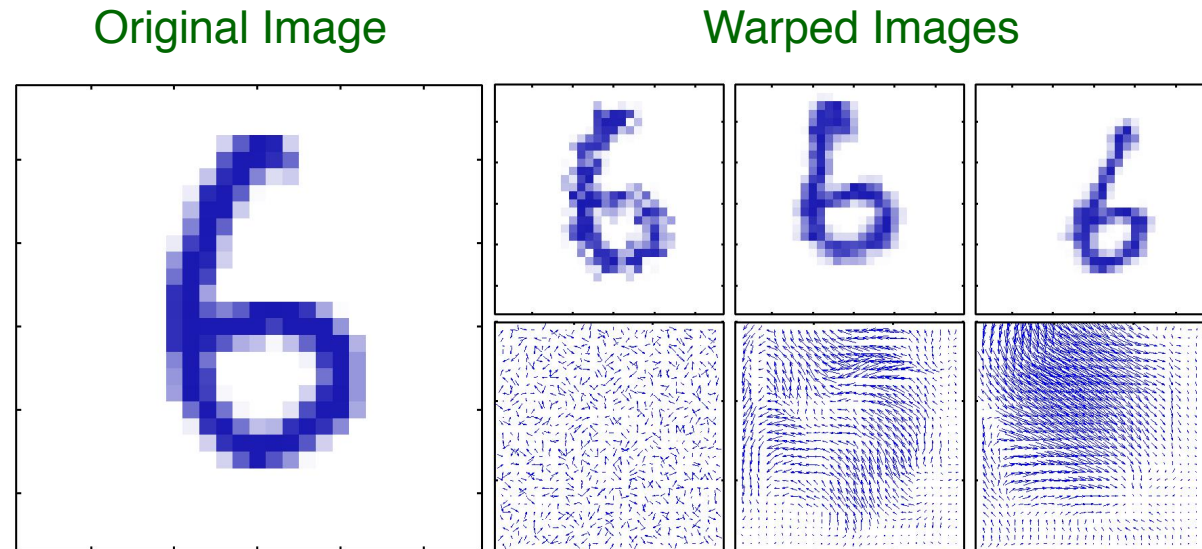
- By extracting features invariant to transformations

4. Build invariance into structure of network (CNN)

- Local receptive fields and shared weights

Approach 1: Transform each input

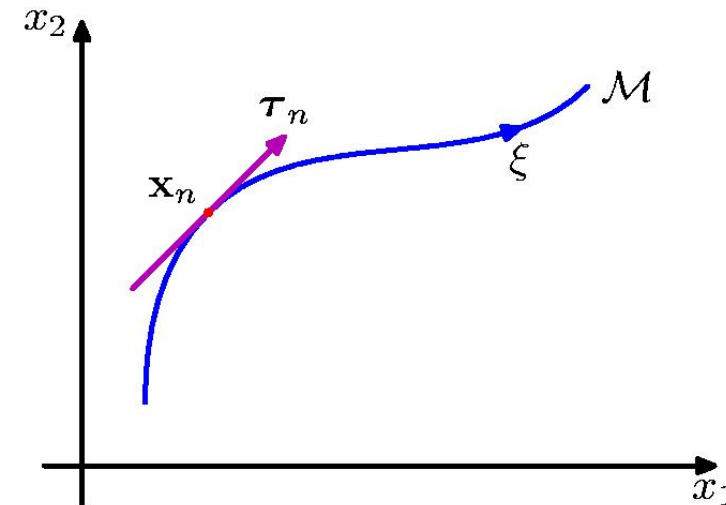
- Synthetically warp each handwritten digit image before presentation to model
- Easy to implement but computationally costly



Random displacements
 $\Delta x, \Delta y \in (0,1)$ at each pixel
Then smooth by convolutions
of width 0.01, 30 and 60 resply

Approach 2: Tangent Propagation

- Regularization to encourage invariance to transformations
 - by technique of tangent propagation
- Consider effect of transformation on input \mathbf{x}_n
 - A 1- D continuous transformation parameterized by ξ applied to \mathbf{x}_n sweeps a manifold \mathcal{M} in D -dimensional input space



Two-dimensional input space showing effect of continuous transformation with single parameter ξ

Let the vector resulting from acting on \mathbf{x}_n by this transformation be denoted by $s(\mathbf{x}_n, \xi)$ defined so that $s(\mathbf{x}, 0) = \mathbf{x}$.

Then the tangent to the curve \mathcal{M} is given by the directional derivative $\tau = \delta s / \delta \xi$ and the tangent vector at point \mathbf{x}_n is given by

$$\tau_n = \left. \frac{\partial s(\mathbf{x}_n, \xi)}{\partial \xi} \right|_{\xi=0}$$

Tangent Propagation as Regularization

- Under a transformation of input vector
 - The network output vector will change
 - Derivative of output k wrt ξ is given by

$$\left. \frac{\partial y_k}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^D \left. \frac{\partial y_k}{\partial x_i} \frac{\partial x_i}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^D J_{ki} \tau_i$$

- where J_{ki} is the (k, i) element of Jacobian Matrix J
- Result used to modify standard error function

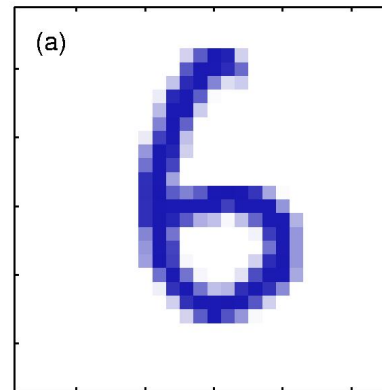
$$\begin{aligned} \tilde{E} &= E + \lambda \Omega \\ \text{where } \lambda &\text{ is a regularization coefficient and} \\ \Omega &= \frac{1}{2} \sum_n \sum_k \left(\left. \frac{\partial y_{nk}}{\partial \xi} \right|_{\xi=0} \right)^2 = \frac{1}{2} \sum_n \sum_k \left(\sum_{i=1}^D J_{nki} \tau_{ni} \right)^2 \end{aligned}$$

- To encourage local invariance in neighborhood of data pt

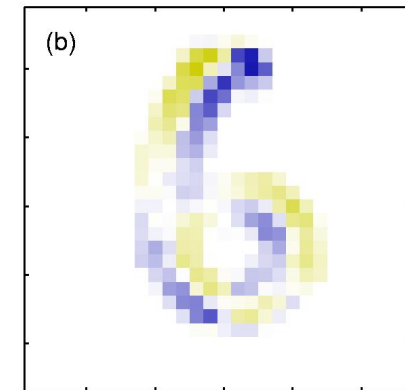
Tangent vector from finite differences

- In practical implementation tangent vector τ_n is approximated using finite differences by subtracting original vector x_n from the corresponding vector after transformations using a small value of ξ and then dividing by ξ
- A related technique *Tangent Distance* is used to build invariance properties into distance-based methods such as *nearest-neighbor* classifiers

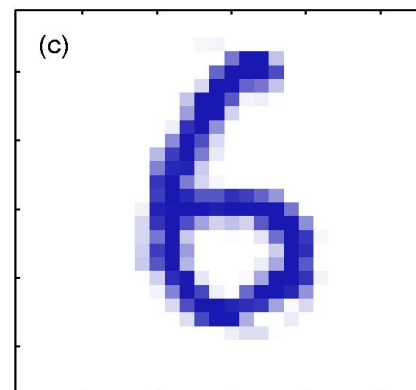
Original Image x



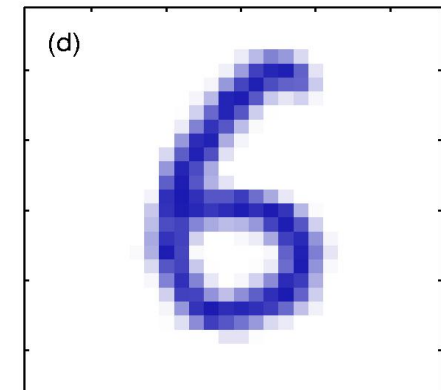
Tangent vector τ corresponding to small clockwise rotation



Adding small contribution from tangent vector to image $x + \epsilon\tau$



True image rotated for comparison



Equivalence of Approaches 1 and 2

- Expanding the training set is closely related to tangent propagation
- Small transformations of the original input vectors together with *Sum-of-squared* error function can be shown to be equivalent to the *tangent propagation* regularizer