Regularization - decrease generalization error, but not training error
Regularization trades increases bias for reduced variance, at a reasonable cost

Limit capacity of a network by adding a parameter norm penalty
$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$ Larger values of alpha = more regularization
We do not induce much variance by leaving the biases unregularized, since it only needs to
observe 1 variable (whereas weights observe how 2 vars interact)
w indicates all the weights, which should be affected, while theta = all the params (including b)

Weight decay / Ridge regression / Tikhonov regularization / L$^2$ parameter norm adds a
regularization term $\frac{1}{2}\|w\|_2^2$ driving weights towards the origin

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^\top w + J(w; X, y),$$

The gradient of the regularized objective function J~ is

with the corresponding parameter gradient $\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y).$

We take a gradient descent step with $w \leftarrow (1 - \epsilon\alpha)w - \epsilon \nabla_w J(w; X, y).$
In a single step, the weight vector is being multiplicatively shrunk

We simplify the objective function by making a quadratic approximation to the objective function
$w^* = \arg\min_w J(w).$

$$\hat{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*)$$

The approximation J^ is given by
H is the Hessian of J w/ respect to w evaluated at w*
Since w* is the minimum where the gradient vanishes, there is no first order term
H is positive semidefinite since w* is the location of a minnimum of J

The minimum of J occurs where its gradient $\nabla_w \hat{J}(w) = H(w - w^*)$ is equal to 0
Modify the above equation by adding the weight decay gradient, and solve for the min

$$\alpha\tilde{w} + H(\tilde{w} - w^*) = 0$$
$$(H + \alpha I)\tilde{w} = Hw^*$$
$$\tilde{w} = (H + \alpha I)^{-1} Hw^*.$$ w-bar represents the location of the minimum
As alpha approaches 0, the regularized solution approaches w*

Decompose H into a diagonal matrix and an orthonormal basis of eigenvectors | $H = Q\Lambda Q^\top.$
$$\tilde{w} = (Q\Lambda Q^\top + \alpha I)^{-1} Q\Lambda Q^\top w^*$$
$$= \left[Q(\Lambda + \alpha I)Q^\top\right]^{-1} Q\Lambda Q^\top w^*$$
$$= Q(\Lambda + \alpha I)^{-1} \Lambda Q^\top w^*.$$

Weight decay rescales w* along the axis defined by eigenvectors of H => the component of w* aligned with the i-th eigenvector of H is rescaled by $\frac{\lambda_i}{\lambda_i + \alpha}$.

When the eigenvalues of H are very large, therefore, lambda$_i$ >> a, regularization does nothing

Only directions along which the parameters contribute significantly to reducing the objective function are preserved; barely significant directions are ignored

For linear regression, the cost function can be expressed in terms of the sum of squared errors $(Xw - y)^\top (Xw - y)$.

With L2 reg, it becomes:
$$(Xw - y)^\top (Xw - y) + \frac{1}{2}\alpha w^\top w.$$

The normal equation for the solution changes from $w = (X^\top X)^{-1} X^\top y$ to
$$w = (X^\top X + \alpha I)^{-1} X^\top y.$$

Without L$^2$, the closed form solution has the covariance matrix $\frac{1}{m} X^\top X$, which is replaced with
$$(X^\top X + \alpha I)^{-1}$$

Diagonal entries of this new matrix correspond to variance; L$^2$ causes the learning algorithm to perceive X as having a higher variance, shrinking weights on features w/ low covariance to the outputs compared to the added variance

L$^1$ regularization is defined as $\Omega(\theta) = ||w||_1 = \sum_i |w_i|,$

$\tilde{J}(w; X, y) = \alpha ||w||_1 + J(w; X, y),$ w/ the gradient

$\nabla_w \tilde{J}(w; X, y) = \alpha \operatorname{sign}(w) + \nabla_w J(X, y; w)$ where sign(w) is the sign of w applied element wise

Regularization contribution does not scale linearly with each w$_i$ but it is a constant factor with a sign equal to sign(w$_i$) => algebraic solutions will not be clean

pg 1

Make a simplifying assumption that the Hessian is diagonal

$H = \operatorname{diag}([H_{1,1}, \ldots, H_{n,n}]),$ where H$_{i,i}$ > 0

This assumption is true if the linreg problem has been preprocessed to remove correlation between input features, which can be done through PCA

The quadratic approximation decomposes into a sum over the params

$$\hat{J}(w; X, y) = J(w^*; X, y) + \sum_i \left[ \frac{1}{2} H_{i,i}(w_i - w_i^*)^2 + \alpha |w_i| \right].$$

$$w_i = \operatorname{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}.$$

The analytic solution is given by

Consider when $w_i^* > 0$ for all i
- If $w_i^* <$ alpha/$H_{i,i}$ => optimal $w_i^*$ is 0 since the contribution of J is overwhelmed by $L_1$ reg, pushing $w_i$ to 0
- If $w_i^* >$ alpha/$H_{i,i}$ regularization shifts it to zero by alpha/$H_{i,i}$

If $w_i^* < 0$, $L^1$ makes $w_i$ less negative by the same shift

$L^1$ results in a more sparse solution
If we add the same constraints for a positive definite, diagonal Hessian to $L^2$ reg, we find that

$$\tilde{w}_i = \frac{H_{i,i}}{H_{i,i}+\alpha} w_i^*$$ If $w_i^*$ is nonzero, then so is the result $\tilde{w}_i$ <= $L^2$ reg does not sparsify

The sparsity propery of $L^1$ is used as a feature selection mechanism; LASSO integrates an $L^1$ penalty with a linear model and least squares cost function, causing a subset of params = 0

Many regularizatoin strategies can be interpreted as MAP Bayesian inference, w/ $L^2$ reg equivalent to MAP with a Gaussian prior on the weights
For $L^1$, the penalty alpha * omega(w) used to regularize a cost function is equivalent to the log prior term maximized by MAP when the prior is an isotropic Laplace distribution over $R^n$

$$\log p(\boldsymbol{w}) = \sum_i \log \text{Laplace}(w_i; 0, \frac{1}{\alpha}) = -\alpha||\boldsymbol{w}||_1 + n\log\alpha - n\log 2.$$

log alpha - log 2 terms do not depend on w and can be ignored

**Norm Penalties as Constrained Optimization**
$$\tilde{J}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) = J(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) + \alpha\Omega(\boldsymbol{\theta}).$$

Recall we can minimize a function subject to constraints by creating a Generalized Lagrange function consisting of the original function plus the penalties, a product between a KKT coefficient, and a function representing whether the constraint is satisfied

To constrain omega(theta) to be less than k: $\mathcal{L}(\boldsymbol{\theta}, \alpha; \boldsymbol{X}, \boldsymbol{y}) = J(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) + \alpha(\Omega(\boldsymbol{\theta}) - k).$

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \alpha).$$
The solution is given by
alpha increases for omega(theta) > k, and decreases if it is less than k; a*, the optimal value, encourages omega(theta) to shrink but not enough to make it less than k
We can fix a* and view the problem as a function of theta

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \alpha^*) = \arg\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) + \alpha^*\Omega(\boldsymbol{\theta}).$$ omega is the norm penalty

A larger alpha corresponds to a smaller constraint region; we do not know the size of the constraint region imposed by alpha since it does not tell us the value of k
To use explicit constraints, we can make grad. descent take a step downhill on J(theta) and project theta back to the nearest point satisfying omega(theta) < k => useful if we do not want to spend time searching for an alpha corresponding to k

Using explicit constraints can prevent units from becoming dead (which penalties do as they can cause the cost function to get stuck in a local minimum)
Explicit constraints also prevent feedback loops (large weights cause large gradients which cause large weights, and so on until numerical overflow)
Constraining the norm of each column of a weight matrix of a neural net layer, rather than the weight matrix's Frobenius norm, makes sure each neuron does not die

## Regularization and Under-Constrained Problems
Many linear models depend on inverting the matrix $X^TX$ which cannot be done if it is singular, whenever the data generating distribution has no variance in a certain direction or when no variance is observed because there are fewer examples than input features

Instead, regularization will focus on inverting $X^\top X + \alpha I$ , which is guaranteed invertible
It is possible for a problem with no closed form solution to be underdetermined
- If linearly separable, logistic regression works for w, 2w, ... , nw and grad. descent will continually increase the weights until overflow <= this is known as underdetermined
- Weight decay causes grad. descent to quit increasing the magnitude of the weights when slope of the likelihood = alpha, the weight decay coefficient

We can solve underdetermined problems using the Moore-Penrose pseudoinverse

$$X^+ = \lim_{\alpha \searrow 0}(X^\top X + \alpha I)^{-1}X^\top$$

This performs linreg with weight decay => it is the limit as the regularization coefficient shrinks to 0, and can be interpreted as stabilizing underdetermined problems using regularization

(Skipped section about dataset augmentation)

## Noise Robustness
Adding noise to the inputs is an effective regularization strategy and is equivalent to imposing a penalty on the norm of the weights, and even beyond that in general
Noise can also be added to weights, which is done in recurrent networks
- Bayesian considers model weights uncertain and representable via a probability distribution which incorporates this uncertainty, which can be done by adding noise

$J = \mathbb{E}_{p(x,y)}\left[(\hat{y}(x) - y)^2\right]$ Noise to weights can also be traditional regularization

m labeled training examples $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}.$

Each input includes a random perturbation $\epsilon_W \sim \mathcal{N}(\epsilon; 0, \eta I)$

The perturbed model is denoted as $\hat{y}_{\epsilon_W}(x).$

$$\tilde{J}_W = \mathbb{E}_{p(x,y,\epsilon_W)} \left[ (\hat{y}_{\epsilon_W}(x) - y)^2 \right]$$
$$= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2(x) - 2y\hat{y}_{\epsilon_W}(x) + y^2 \right].$$

The objective function becomes

For a small learning rate, the minimization of J with added weight noise is equivalent to

minimization of J with an additional regularization term: $\eta \mathbb{E}_{p(x,y)} \left[ \| \nabla_W \hat{y}(x) \|^2 \right]$

This encourages parameters to go to regions where small weight perturbations have a small influence on the output => minima surrounded by flat regions

For linreg, this simplifies to $\eta \mathbb{E}_{p(x)} \left[ \| x \|^2 \right]$, which is not a function of params and does not contribute to the gradient of J~

Label smoothing regularizes a model based on softmax w/ k output values by replacing the 0

and 1 classification targets with $\frac{\epsilon}{k-1}$ and 1 - eps. since there is a small chance each label is mistaken
Label smoothing prevents the pursuit of hard probabilities, which are impossible to get

**Semi-Supervised Learning**
Semi-supervised learning uses unsupervised and supervised examples to predict P(y|x)

In deep learning it refers to learning a representation $h = f(x)$.
Learns a representation such that examples in the same class have similar representations
Can construct generative models of P(x) or P(x, y) which shares parameters with a discriminative model of P(y|x)
The supervised criterion can be traded off w/ the unsupervised one
   ● Assumes that P(x) is connected to structure of P(y|x) by shared parameters

**Multi-Task Learning**
Multi-task learning pools examples arising out of similar tasks to constrain the parameters and generalize better
Different tasks y which take the same input x can have the same representation h
Task specific parameters are for a specific task, and are the upper layers of a neural net
Generic parameters, shared across similar tasks, are the lower layers
The statistical strength of the shared parameteers is due to extra training examples
The prior: among factors explaining variations observed in data associated with different tasks, some are shared across two or more tasks

**Early Stopping**
Return to the point with the lowest validation error, rather than training error to prevent overfit

Save the parameters every time the model improves its validation error

Let $n$ be the number of steps between evaluations.
Let $p$ be the "patience," the number of times to observe worsening validation set error before giving up.
Let $\theta_o$ be the initial parameters.
$\theta \leftarrow \theta_o$
$i \leftarrow 0$
$j \leftarrow 0$
$v \leftarrow \infty$
$\theta^* \leftarrow \theta$
$i^* \leftarrow i$
**while** $j < p$ **do**
   Update $\theta$ by running the training algorithm for $n$ steps.
   $i \leftarrow i + n$
   $v' \leftarrow \text{ValidationSetError}(\theta)$
   **if** $v' < v$ **then**
      $j \leftarrow 0$
      $\theta^* \leftarrow \theta$
      $i^* \leftarrow i$
      $v \leftarrow v'$
   **else**
      $j \leftarrow j + 1$
   **end if**
**end while**
Best parameters are $\theta^*$, best number of training steps is $i^*$

The number of training steps can be considered another hyperparameter
Ideally, the validation set is being processed on a different GPU in parallel to improve efficiency
Once early stopping is used to determine the optimal # of steps, the validation data can be used for training as well

**Algorithm 7.2** A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

---

Let $X^{(\text{train})}$ and $y^{(\text{train})}$ be the training set.
Split $X^{(\text{train})}$ and $y^{(\text{train})}$ into $(X^{(\text{subtrain})}, X^{(\text{valid})})$ and $(y^{(\text{subtrain})}, y^{(\text{valid})})$ respectively.
Run early stopping (algorithm 7.1) starting from random $\theta$ using $X^{(\text{subtrain})}$ and $y^{(\text{subtrain})}$ for training data and $X^{(\text{valid})}$ and $y^{(\text{valid})}$ for validation data. This returns $i^*$, the optimal number of steps.
Set $\theta$ to random values again.
Train on $X^{(\text{train})}$ and $y^{(\text{train})}$ for $i^*$ steps.

Early stopping is a good regularizer as it does not require its gradient to be computed
Early stopping acts as a regularizer by restricting the optimization procedure to a small volume of parameter space around the initial weights $\theta_0$
- The effective capacity is given by eT, where e is the learning rate and T is the # of steps
- eT behaves as if it were the reciprocal of the weight decay coefficient

Consider a linear model, where the only params are weights (not biases)

The cost function J(Θ) can be modeled w/ a quadratic approximation near w*, the optimal wts

$$\hat{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*),$$

If w* is a minimum of J(w), the hessian H is positive semidefinite
A taylor series approximation gives the gradient as:

$$\nabla_w \hat{J}(w) = H(w - w^*).$$

Let the initial parameter vector be **0**
Gradient descent does the following
pg 2

$$w^{(\tau)} = w^{(\tau-1)} - \epsilon \nabla_w \hat{J}(w^{(\tau-1)})$$
$$= w^{(\tau-1)} - \epsilon H(w^{(\tau-1)} - w^*)$$
$$w^{(\tau)} - w^* = (I - \epsilon H)(w^{(\tau-1)} - w^*).$$

Rewrite this expression in the space of eigenvectors of H, using the fact $H = Q\Lambda Q^\top$ where Q is an orthonormal basis of eigenvectors

$$w^{(\tau)} - w^* = (I - \epsilon Q\Lambda Q^\top)(w^{(\tau-1)} - w^*)$$
$$Q^\top(w^{(\tau)} - w^*) = (I - \epsilon\Lambda)Q^\top(w^{(\tau-1)} - w^*)$$

Assume that $w^0$ is **0** and epsilon is small enough to guarantee $|1 - \epsilon\lambda_i| < 1$ (which allows us to say [I - [...])$^\top$], the parameter trajectory during training is as follows

$$Q^\top w^{(\tau)} = [I - (I - \epsilon\Lambda)^\tau]Q^\top w^*.$$

Rearrange the equation for Q$^\top$w~ from L$^2$ regularization

$$Q^\top \tilde{w} = (\Lambda + \alpha I)^{-1}\Lambda Q^\top w^*$$
$$Q^\top \tilde{w} = [I - (\Lambda + \alpha I)^{-1}\alpha]Q^\top w^*$$

We can see that if the hyperparameters epsilon, T, and alpha are chosen such that

$$(I - \epsilon\Lambda)^\tau = (\Lambda + \alpha I)^{-1}\alpha,$$ , then early stopping is equivalent to L$^2$ reg

If all lambda$_i$ are small, $\epsilon\lambda_i \ll 1 \text{ and } \lambda_i/\alpha \ll 1$

Number of training iterations is inversely proportional to L$^2$ reg param, and the inverse of T(epsilon) plays the weight decay coefficient
Relating early stopping to L$^2$ reg, values corresponding to directions of significant curvature learn early relative to directions of less curvature

**Parameter Tying and Parameter Sharing**

We can assert a prior by imposing dependencies between model parameters
Suppose we have two models performing the same classification task, with different input distributions: model A w/ parameters $w^A$ and model B w/ parameters $w^B$

The two models have related outputs as follows: $\hat{y}^{(A)} = f(w^{(A)}, x)$ ; $\hat{y}^{(B)} = g(w^{(B)}, x)$.
If the tasks are similar enough, we say that for all i, $w_i^A = w_i^B$

We use a parameter $L^2$ norm penalty: $\Omega(w^{(A)}, w^{(B)}) = \|w^{(A)} - w^{(B)}\|_2^2.$

Another way to regularize params to be close to one another is to use constraints: to force sets of parameters to be equal
This regularization method is known as parameter sharing, since various models share a unique set of parameters
Only a subset of parameters, the unique set, for each model needs to be stored in memory

CNNs share parameters across multiple image locations, since photos are invariant to translation => the feature (hidden unit w/ same weights) is computed over various locations
This significantly lowers # of unique model params and significantly increases network sizes w/o requiring an increase in training data

**Sparse Representations**
Placing a penalty on the activations of a network, encouraging them to be sparse imposes a complicated penalty on the model's parameters
$L^1$ reg induces a sparse parametrization, encouraging many of the params to be 0
Representational sparsity describes a representation where many of the elements are 0 or close to 0; the below is a sparsely parametrized linear model

$$
\begin{bmatrix} 18 \\ 5 \\ 15 \\ -9 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 & 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & -1 & 0 & 3 & 0 \\ 0 & 5 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & -4 \\ 1 & 0 & 0 & 0 & -5 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ -2 \\ -5 \\ 1 \\ 4 \end{bmatrix}
$$
$$
y \in \mathbb{R}^m \qquad A \in \mathbb{R}^{m \times n} \qquad x \in \mathbb{R}^n
$$

The one below is linear regression w/ a sparse representation h of the data x

$$
\begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix}
$$
$$
y \in \mathbb{R}^m \qquad B \in \mathbb{R}^{m \times n} \qquad h \in \mathbb{R}^n
$$

h represents the info in x, but does so with a sparse vector

Adding to the loss function J a norm penalty on the representation can achieve representational sparsity: $\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha\Omega(h)$

$\Omega(h) = ||h||_1 = \sum_i |h_i|.$ $L^1$ is only one choice of penalty which can result in sparse rep.

Hard constraint on activations can achieve representational sparsity
Orthogonal matching pursuit encodes an input x with the representation h which solves the constrained optimization problem

$$\underset{h, ||h||_0 < k}{\arg\min} ||x - Wh||^2.$$

where $||h||_0$ is the # of nonzero entries of h
This problem can be efficiently solved when W is orthogonal; this method is called OMP-k w/ k indicating the # of nonzero features allowed

**Bagging and Other Ensemble Methods**
Bagging, or bootstrap aggregating, reduces generalization error by combining several methods which then vote on the output for the test set
Averaging works since different models won't make the same errors on the test set

Consider a set of k regression models, each making an error of $\epsilon_i$ with the errors drawn from a normal distribution with variance $\mathbb{E}[\epsilon_i^2] = v$ and covariance $\mathbb{E}[\epsilon_i \epsilon_j] = $
The error made by the average prediction of all ensemble methods is 1/k($\Sigma_i$ epsilon$_i$)
The expected squared error is given by

$$\mathbb{E}\left[\left(\frac{1}{k}\sum_i \epsilon_i\right)^2\right] = \frac{1}{k^2}\mathbb{E}\left[\sum_i \left(\epsilon_i^2 + \sum_{j \neq i} \epsilon_i \epsilon_j\right)\right]$$
$$= \frac{1}{k}v + \frac{k-1}{k}c.$$

When the errors are perfectly correlated, and c = v, the expected squared error is v, so model averaging does not help at all
When the errors are perfectly uncorrelated and c = 0, the expected squared error is v/k; the ensemble will perform at least as well as any of its members, and much better if the errors made are independent

Bagging allows the same kind of model, training objective, and objective function to be used several times
Bagging involves constructing k different datasets, each with the same # of examples as the original dataset, sampled with replacement; each dataset is missing some examples from the original dataset, and contains multiple duplicates, causing variance between each model
Even networks trained on the same examples will have multiple solution points as a result of random initialization, causing different members to make partially independent errors
Benchmark comparisons are made using a single model, not ensembles

Boosting constructs an ensemble w/ a higher capacity than each individual model, building ensembles of neural nets by adding networks to the ensemble

**Dropout**
Dropout provides an inexpensive approximation to training and lets us evaluate a bagged ensemble of exponentially many networks
Dropout trains the ensemble consisting of all networks that can be formed by removing non-output units from an underlying base network
We can remove a unit from a network by multiplying its output by 0 (in most cases, but this may need to be trivially modified when dealing with things like RBF)

Dropout aims to replicate bagging; using mini-batch learning algorithm, we randomly sample a different binary mask to apply to the input and hidden units
The probablity of sampling a mask value of 1, indicating the unit should be included, is a hyperparameter; also the probability of whether an input is used or not can be adjusted
We then run forward prop, and backprop as usual

Suppose a mask vector $\mu$ specifies which units to include, and $J(\Theta, \mu)$ defines the cost of the model defined by parameters $\Theta$ and mask $\mu$
Dropout training consists of minimizing $E_\mu J(\Theta, \mu)$; this expectation contains exponentially many terms (all possible masks) but we can obtain an unbiased estimate by sampling values of $\mu$

In dropout, each parent neural network inherits a different subset of params from its parent neural network, making it possible to represent an exponential number of models
With dropout, a tiny fraction of the sub-networks are each trained for a single step, and parameter sharing causes the remaining sub-networks to arrive at good parameter values
Beyond these differences, dropout is essentially the same as bagging

A bagged ensemble must accumulate votes from all its models - this is known as inference
In the case of bagging, each model produces a probability distribution $p^i(y \mid x)$; the prediction of the ensemble is given by the arithmetic mean of these distributions

$$\frac{1}{k} \sum_{i=1}^{k} p^{(i)}(y \mid x).$$

W/ dropout, each sub-model defined by a mask $\mu$ defines the probability

distribution $p(y \mid x, \mu)$; the arithmetic mean over all masks is given by $\sum_{\mu} p(\mu) p(y \mid x, \mu)$ where $p(\mu)$ is the probability distribution used to sample $\mu$ at training time
We can approximate the inference with sampling, by averaging the output from 10-20 masks
However, it is possible to obtain a good approximation in one forward prop, by using the geometric mean rather than the arithmetic mean of the ensemble members' predictions

The geometric mean of multiple prob. distributions is not guaranteed to be a probability distribution, so we impose the requirement that the prediction of any event cannot = 0, and we renormalize the resulting product distribution

The unnormalized prob. distribution defined by the geometric mean is given by

$$\tilde{P}_{ensemble}(y \mid x) = \sqrt[2^d]{\prod_{\mu} p(y \mid x, \mu)}$$

where d is the # of units that may have been dropped

TO make predictions, the ensemble must be renormalized

$$P_{ensemble}(y \mid x) = \frac{\tilde{P}_{ensemble}(y \mid x)}{\sum_{y'} \tilde{P}_{ensemble}(y' \mid x)}.$$

We can approximate $p_{ensemble}$ by evaluating p(y|x) in one model, which contains all units but with the weights going out unit i multiplied by the probability of including unit i

This captures the right expected value of the unit's output - this is known as the weight scaling inference rule, which performs well empirically

Since dropout is usually ½, the weight scaling rule results in dividing the weights by 2 at the end of training, which makes sure that the total expected input at test time (when all units are used) is roughly the same as train time

Consider a softmax regression w/ n inputs represented by vector v:

$$P(y = y \mid \mathbf{v}) = \text{softmax}\left(\mathbf{W}^\top \mathbf{v} + \mathbf{b}\right)_y.$$

Multiply with a binary vector d (element wise) to index into the family of the sub-models

$$P(y = y \mid \mathbf{v}; \mathbf{d}) = \text{softmax}\left(\mathbf{W}^\top (\mathbf{d} \odot \mathbf{v}) + \mathbf{b}\right)_y.$$

Ensemble predictor defined by renormalizing geometric mean over all ensemble predictions

$$P_{ensemble}(y = y \mid \mathbf{v}) = \frac{\tilde{P}_{ensemble}(y = y \mid \mathbf{v})}{\sum_{y'} \tilde{P}_{ensemble}(y = y' \mid \mathbf{v})}$$

Denominator represents all predictions

$$\tilde{P}_{ensemble}(y = y \mid \mathbf{v}) = \sqrt[2^n]{\prod_{d \in \{0,1\}^n} P(y = y \mid \mathbf{v}; \mathbf{d})}.$$

where

Simplify P~$_{ensemble}$ to see the weight scaling rule is exact

$$\tilde{P}_{ensemble}(y = y \mid \mathbf{v}) = \sqrt[2^n]{\prod_{d \in \{0,1\}^n} P(y = y \mid \mathbf{v}; \mathbf{d})}$$

$$= \sqrt[2^n]{\prod_{d \in \{0,1\}^n} \text{softmax}\left(\mathbf{W}^\top (\mathbf{d} \odot \mathbf{v}) + \mathbf{b}\right)_y}$$

$$= \sqrt[2^n]{\prod_{d \in \{0,1\}^n} \frac{\exp\left(\mathbf{W}_{y,:}^\top (\mathbf{d} \odot \mathbf{v}) + b_y\right)}{\sum_{y'} \exp\left(\mathbf{W}_{y',:}^\top (\mathbf{d} \odot \mathbf{v}) + b_{y'}\right)}}$$

$$= \frac{\sqrt[2^n]{\prod_{d \in \{0,1\}^n} \exp\left(\mathbf{W}_{y,:}^\top (\mathbf{d} \odot \mathbf{v}) + b_y\right)}}{\sqrt[2^n]{\prod_{d \in \{0,1\}^n} \sum_{y'} \exp\left(\mathbf{W}_{y',:}^\top (\mathbf{d} \odot \mathbf{v}) + b_{y'}\right)}}$$

Since P˜ will be renormalized, we can ignore multiplication by factors that are constant in y

$$\tilde{P}_{\text{ensemble}}(y = y \mid \mathbf{v}) \propto \sqrt[2^n]{\prod_{\mathbf{d} \in \{0,1\}^n} \exp\left(W_{y,:}^\top (\mathbf{d} \odot \mathbf{v}) + b_y\right)}$$

$$= \exp\left(\frac{1}{2^n} \sum_{\mathbf{d} \in \{0,1\}^n} W_{y,:}^\top (\mathbf{d} \odot \mathbf{v}) + b_y\right)$$

$$= \exp\left(\frac{1}{2} W_{y,:}^\top \mathbf{v} + b_y\right).$$

=> Discard the denominator

Substituting this into the initial renormalized equation, the softmax classifier has weights ½W

The weight scaling rule is precise in other settings like regression w/ conditionally normal outputs and deep networks without nonlinearities
The approximation is empirically effective for deep nets with nonlinearities
Weight scaling can work better than Monte Carlo approximations, with just 20 samples

Dropout is computationally cheap - $O(n)$ to multiply the binary mask and $O(n)$ to store it, if the implementation requires it; weights must be divided by 2 to start running inference
Dropout does not limit model choice, as it works with neural nets, Boltzmann machines, etc.
However, dropout requires a larger model in order to obtain marginal boosts in validation sets
Unsupervised learning and Bayesian networks can be better than dropout if only few samples

When applied to linreg, dropout is the same as weight decay, with alpha being determined by each input's variance
Stochasticity while training w/ dropout is not necessary (it just approximates the sum over all submodels); fast dropout, reducing stochasticity, is a more principled approximation which can match standard dropout for small datasets

Training the network as a single model with dropout is better than dropout boosting (which is analogous to boosting, with multiple models and little regularization effect)
Regularization effect of bagged ensemble is effective when stochastically sampled members are trained to perform well independently of one another

DropConnect - each product between a single weight and hidden unit can be dropped
Stochastic pooling - ensembles of CNNs w/ each one attending to different spatial locations
Training a network w/ stochastic behavior and averaging over multiple decisions implements bagging with parameter sharing

Dropout is bagging an ensemble of models formed by any random modification, which neural networks can learn to resist (and ideally, those which allow quick approximate inference)
Any modification parametrized by $\mu$ is training an ensemble consisting of $p(y \mid x, \mu)$ for every $\mu$

$\mu \sim N(1, I)$ can outperform dropout on binary masks, since $E(\mu) = 1$, approximate wt. scaling is automatically implemented

Furthermore, dropout trains an ensemble of models that share hidden units, meaning they must be able to perform well regardless of which other hidden units are used in the model, causing the hidden units to be features which are applicable in many contexts

Dropout is highly destructive to the information content of the inputs, rather than the raw inputs themselves; for example, if a model learns an activation which encodes a nose, dropping $h_i$ corresponds to dropping a nose, forcing a model to identify a face using some other feature
This is different than just dropping some features of the nose, as traditional noise does

Dropout uses multiplicative noise, rather than additive, preventing the model from causing $h_i$ to become very large in order to make the noise epsilon insignificant
Batch normalization reparametrizes the model w/ multiplicative and additive noise to improve optimization; however, the noise can have a regularizing effect making dropout unnecessary

**Adversarial Training**
Even neural nets performing at human accuracy have a nearly 100% error rate on examples intentionally constructed using an optimization procedure searching for an input x' near an arbitrary point x such that the model output is very different at x'
Many times, x' can be so similar to x that a human cannot tell the difference, but the network can distinguish between the original and adversarial network

One can reduce the error rate on the i.i.d test set through adversarial training, which trains on adversarially perturbed examples from the training set
One of the cuases of adversarial examples is excessive linearity
Neural nets are primarily linear, which makes them easy to optimize; however, the value of a linear function changes rapidly by epsilon*$||w||_1$ where epsilon is a small change in inputs
This can compound and cause the model to predict something entirely different
Adversarial training encourages the network to be locally constant in the area of training data

Purely linear models like logistic regression cannot resist adversarial examples as they are forced to be linear; the flexibility of neural nts to learn linear / nonlinear functions allows them to capture linear trends in training data while being able to resist local perturbation

Adversarial examples generated using a label provided by a trained model are called virtual adversarial examples; the classifier may then be trained to assign the same label to x and x'
Requires assumption that small changes in the manifold do not change the class to change

**Tangent Distance, Tangent Prop, and Manifold Tangent Classifier**
Recall manifold hypothesis - all data lies near a low dimensional manifold
Tangent distance - metric derived from knowledge of manifolds where probability concentrates

Nearest neighbor distance between $x_1$ and $x_2$ corresponds to distance between $M_1$ and $M_2$ which is their respective manifolds

A cheap alternative to this is to approximate $M_i$ by the tangent plane at $x_i$ and measure the distance between the two tangents or a tangent plane and a point, which can be achieved by solving a low dimensional (same as manifolds) linear system

Tangent prop trains a classifier with an extra penalty to make each output locally invariant to known factors of variation

These factors of variation correspond to movement along the manifold near which examples of the same class concentrate; local invariance requires $\nabla_x f(x)$ to be orthogonal to the known manifold vectors $v^i$ at x, or the directional derivative of f at x in the directions $v^i$ be small, which can be accomplished by adding a small regularization penalty omega

$$\Omega(f) = \sum_i \left( \left( \nabla_x f(x) \right)^\top v^{(i)} \right)^a$$

Tangent vectors are derived from the knowledge of effects of transformations such as rotations

Tangent prop is also useful in the context of supervised learning

In dataset augmentation, the network is explicitly trained to correctly classify distinct inputs created by applying more than an infinitesimal amount of these transformations

Tangent prop regularizes the model to resist perturbation in directions correspoding to the transformations, which only regularizes the model to resist infinitesimal perturbation

The infinitesimal approach can only shrink derivatives of ReLUs by turning units off or shrinking weights, but not by saturating at a high value like how tanh works

Dataset augmentation works well with ReLUs since different subsets of rectified units activate for different transformed versions of each input

Double backprop regularizes the Jacobian to be small

Double backprop and adversarial training require the model should be invariant to all directions of change, while dataset augmentation and tangent prop require model to be invariant to certain specified directions

Adversarial training is the non-infinitesimal version of double backprop

Manifold tangent classifier eliminates the need to know the tangent vectors before; it makes use of autoencoders, which can be used to estimate the manifold tangent vectors

These estimated tangent vectors go beyond classical invariants such as rotation

Manifold tangent classifier uses an autoencoder to learn the manifold structure w/ unsupervised learning, and uses the tangents to regularize a neural net as in tangent prop