
MACHINE LEARNING AND PATTERN RECOGNITION

Fall 2005, Lecture 2:

Energy-Based Models and Loss Functions,
Linear Machines

Yann LeCun
The Courant Institute,
New York University
<http://yann.lecun.com>

Linear Machines: Regression with Mean Square

Linear Regression, Mean Square Loss:

- decision rule: $y = W'X$
- loss function: $L(W, y^i, X^i) = \frac{1}{2}(y^i - W'X^i)^2$
- gradient of loss: $\frac{\partial L(W, y^i, X^i)}{\partial W} = -(y^i - W(t)'X^i)X^i$
- update rule: $W(t+1) = W(t) + \eta(t)(y^i - W(t)'X^i)X^i$
- direct solution: solve linear system $[\sum_{i=1}^P X^i X^{i'}]W = \sum_{i=1}^P y^i X^i$

Linear Machines: Perceptron

Perceptron:

- decision rule: $y = F(W'X)$ (F is the threshold function)
- loss function: $L(W, y^i, X^i) = (F(W'X^i) - y^i)W'X^i$
- gradient of loss: $\frac{\partial L(W, y^i, X^i)}{\partial W} = -(y^i - F(W(t)'X^i))X^i$
- update rule: $W(t+1) = W(t) + \eta(t)(y^i - F(W(t)'X^i))X^i$
- direct solution: find W such that $-y^i F(W'X^i) < 0 \quad \forall i$

Linear Machines: Logistic Regression

Logistic Regression, Negative Log-Likelihood Loss function:

- decision rule: $y = F(W'X)$, with $F(a) = \frac{1 - \exp(a)}{1 + \exp(a)}$ (sigmoid function).
- loss function: $L(W, y^i, X^i) = 2 \log(1 + \exp(-y^i W'X^i))$
- gradient of loss: $\frac{\partial L(W, y^i, X^i)}{\partial W} = -(Y^i - F(W'X)) X^i$
- update rule: $W(t + 1) = W(t) + \eta(t)(y^i - F(W(t)'X^i))X^i$

General Gradient-Based Supervised Learning Machine

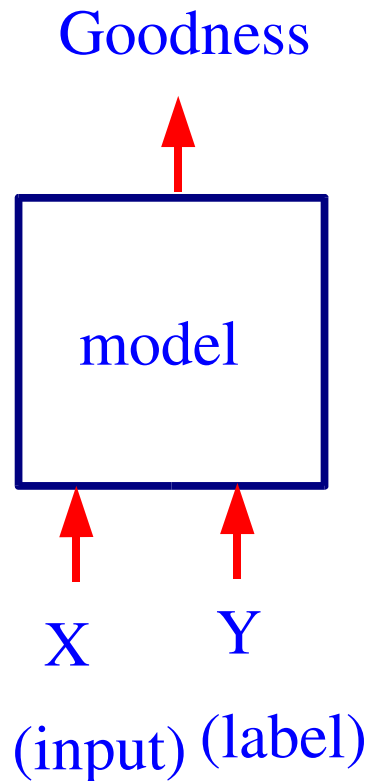
Neural Nets, and many other models:

- decision rule: $y = F(W, X)$, where F is some function, and W some parameter vector.
- loss function: $L(W, y^i, X^i) = D(y^i, F(W, X^i))$, where $D(y, f)$ measures the “discrepancy” between A and B .
- gradient of loss: $\frac{\partial L(W, y^i, X^i)}{\partial W} = \frac{\partial D(y^i, f)}{\partial f} \frac{\partial F(W, X^i)}{\partial W}$
- update rule: $W(t + 1) = W(t) - \eta(t) \frac{\partial D(y^i, f)}{\partial f} \frac{\partial F(W, X^i)}{\partial W}$

Three Questions:

- What architecture $F(W, X)$.
- What loss Function $L(W, y^i, X^i)$.
- What optimization method.

A Model is Designed and trained to Answer Questions



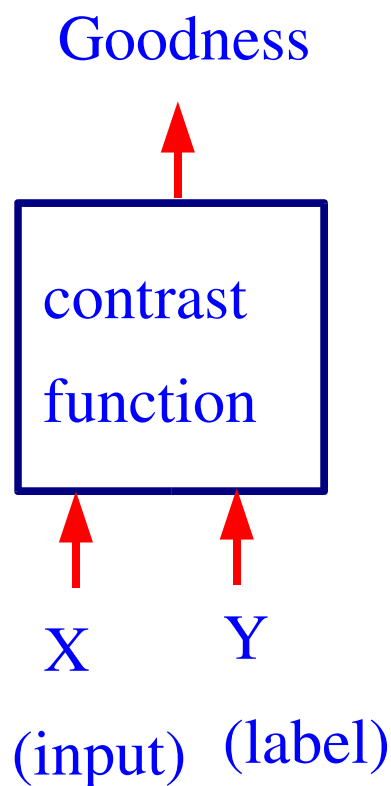
- **Example:** X is an image from a camera; Y is a discrete variable e.g. Y in {animal, human, plane, truck, car}.
- **1. Best Guess for Y:** which category best describes X?
- **2. Ranking on Y:** Is X more a car than an airplane?
- **3. Distribution on Y:** give an estimate of $P(\text{animal} | X)$
- **4. Best Guess for X:** among all images of airplane, give me the best one.
- **5. Ranking on X:** is this image more of a truck than that one?
- **6. Distribution on X:** among all images of airplanes, how likely is this one?

For each question, a different learning strategy is required

Do not answer a more complex question than necessary

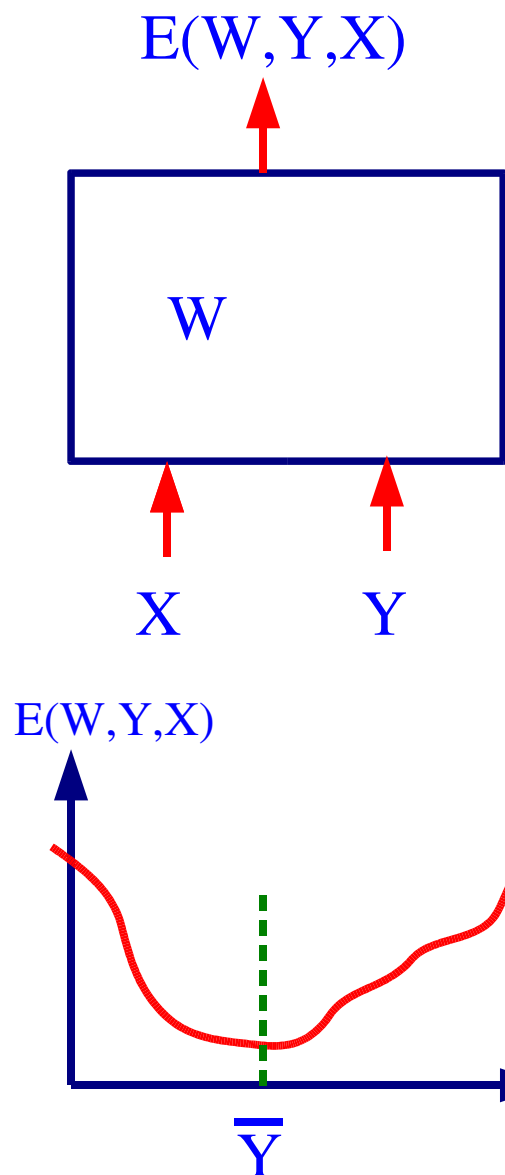
1. **Best Guess for Y:** which category best describes X?
 2. **Ranking on Y:** Is X more a car than an airplane?
 3. **Distribution on Y:** give an estimate of $P(\text{animal} | X)$
 4. **Best Guess for X:** among all images of airplane, give me the best one.
 5. **Ranking on X:** is this image more of a truck than that one?
 6. **Distribution on X:** among all images of airplanes, how likely is this one?
- The questions are in increasing order of complexity.
- The machine should be designed and trained to answer the simplest question possible.

What is a model?



- A model measures the **goodness** of a combination of **observed variables (X)** and **variables to be predicted (Y)**.
- Probabilistic approaches compute the **distribution $P(Y|X)$** , and choose the Y that maximizes it.
- Sometimes, we do not need probabilities.
- **Example:** driving a robot. When the robot faces an obstacle, it **MUST** turn left of right. Computing a distribution of steering angles is of little use.
- **Question:** why estimate the whole distribution $P(Y|X)$ when we are only interested in picking the best value of Y?

Energy-Based Models



- **$E(W, Y, X)$** : is a scalar **energy function** (a.k.a. Contrast function) that measures the “compatibility” between Y and X .
- **W** is the parameter to be learned.
- **MAP Inference**: Given an input X , find the value of Y that minimizes the energy:

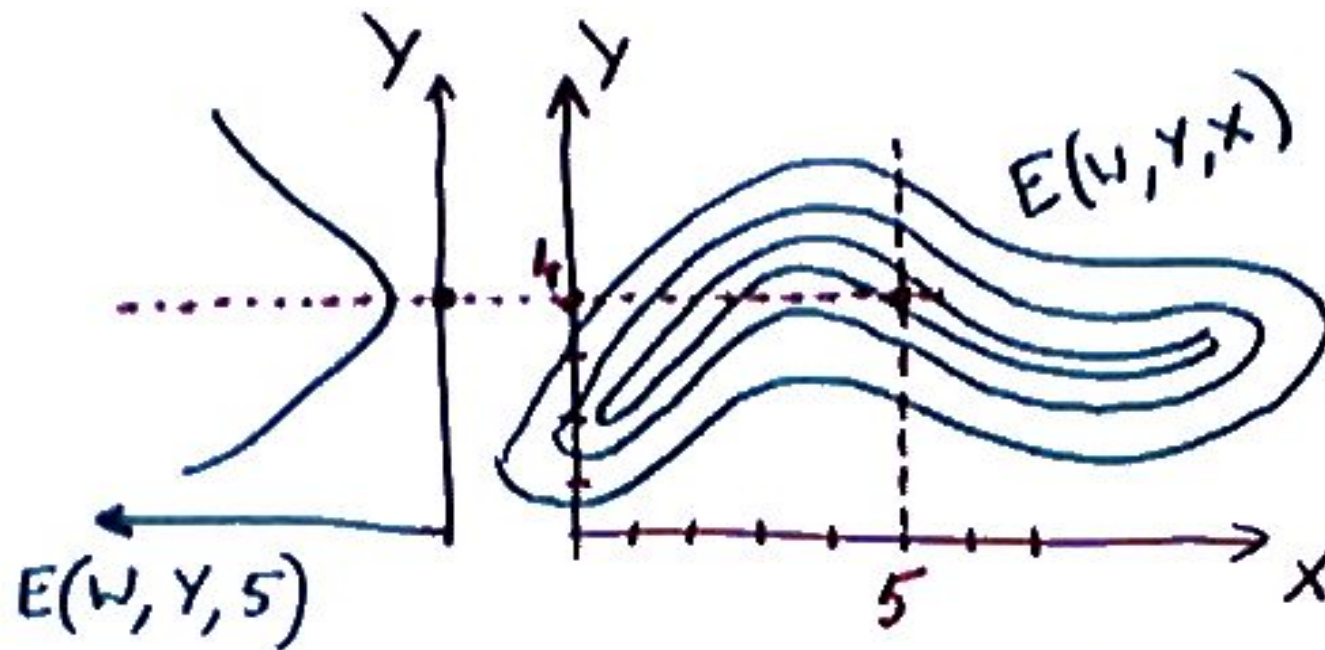
$$\check{Y} = \operatorname{argmin}_{y \in \{Y\}} E(W, y, X)$$

- **Probabilistic Prediction**: Given an input X , compute the conditional distribution over Y (Gibbs Distribution):

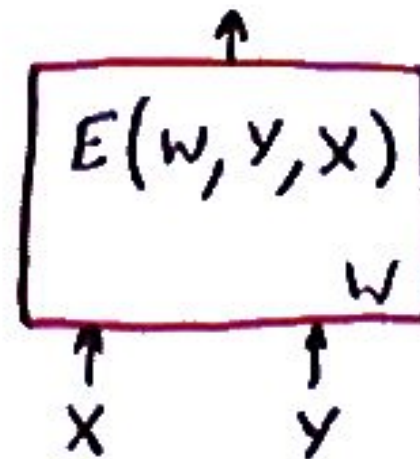
$$P(Y|X) = \frac{\exp(-\beta E(W, Y, X))}{\sum_{y \in \{Y\}} \exp(-\beta E(W, y, X))}$$

- **For decision making, we need no normalization.**

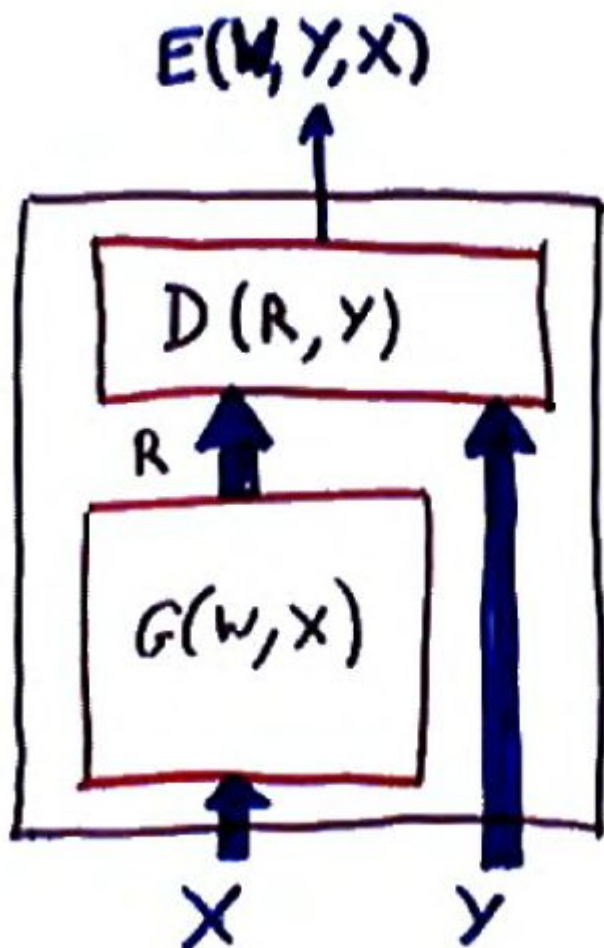
MAP Inference with Energy-Based Models



$$\hat{y} = \min_y E(w, y, x)$$

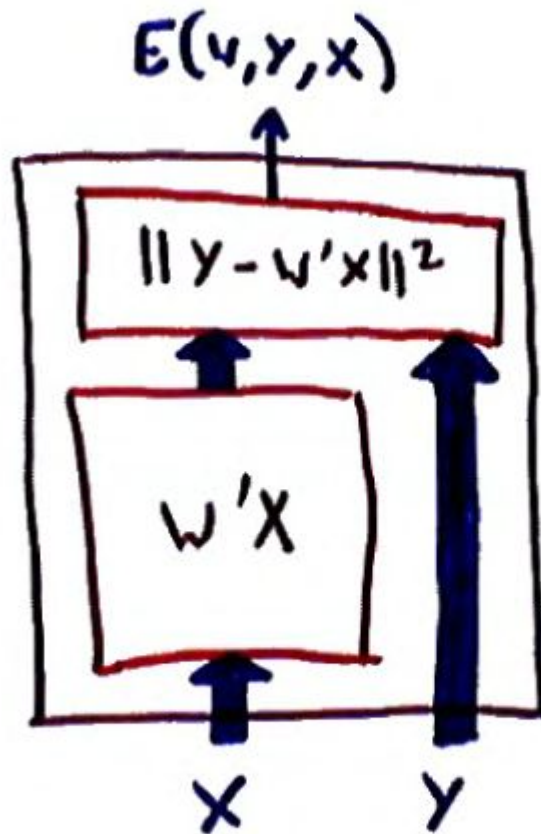


Examples of EBM: Regressor



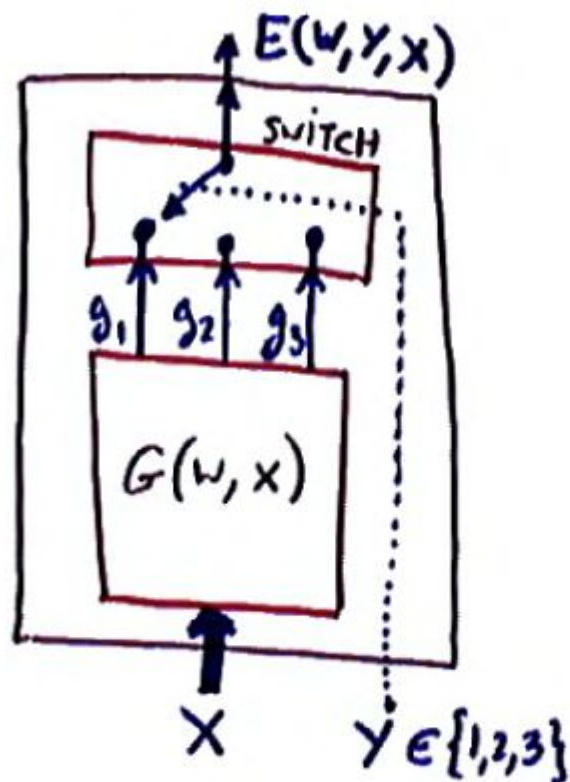
- X and Y are vectors or other entities
- Energy: $E(W, Y, X) = D(Y, G(W, X))$ where $D(Y, R)$ is a distance or dissimilarity measure.
- Best output: $\check{Y} = \min_Y E(W, Y, X) = G(W, X)$.

Examples of EBM Regressor: Linear Regression



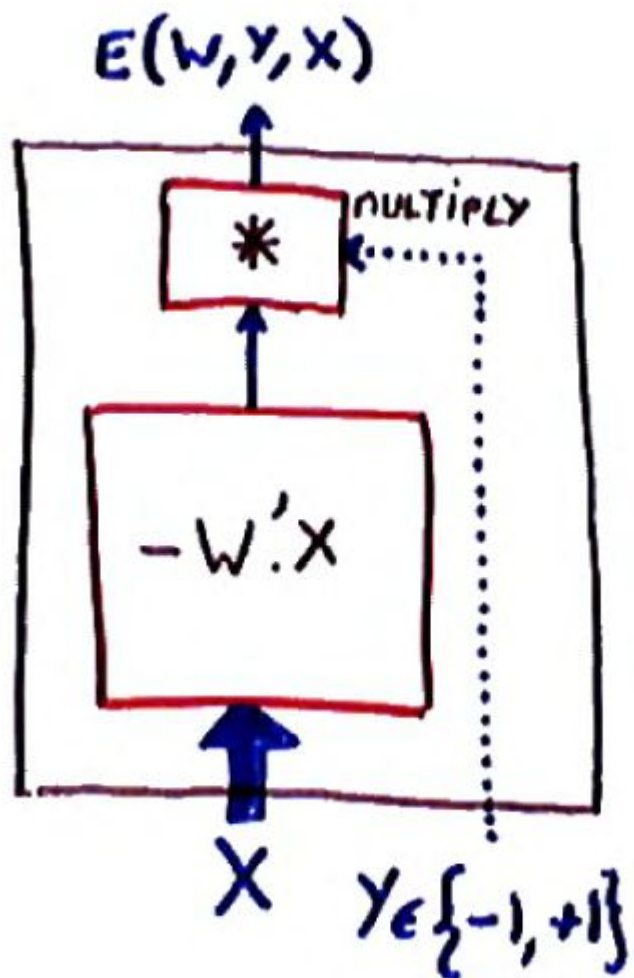
- X and Y are vectors
- Energy: $E(W, Y, X) = \|Y - W'X\|^2$.
- Best output: $\check{Y} = \min_Y E(W, Y, X) = W'X$.

Examples of EBM: Classifier



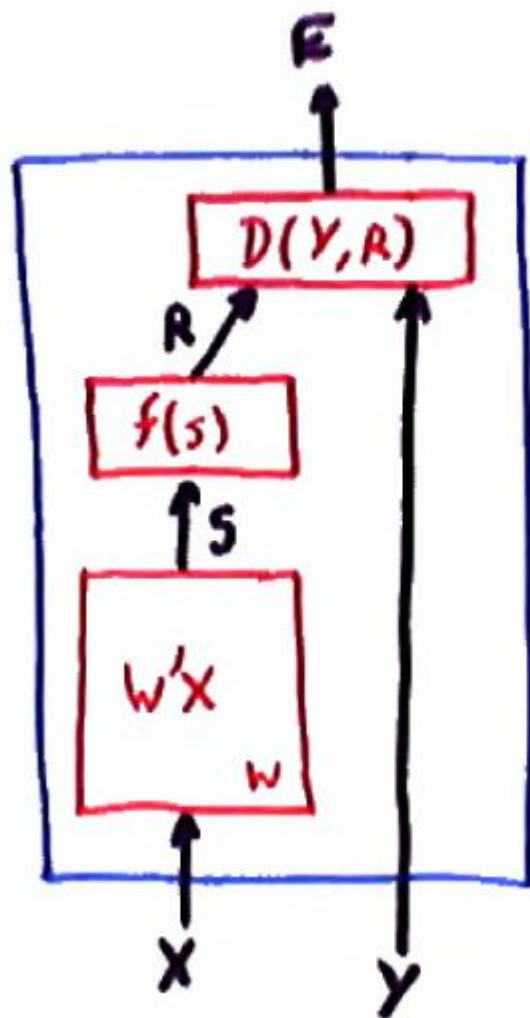
- Y is a discrete variable, $\{Y\} = \{1, 2, 3\}$.
- Energy: $E(W, Y, X) = \sum_k G_k(W, X) \delta(k, Y)$, where $\delta(k, Y) = 1$ iff $k = Y$ and 0 otherwise.
- $G_k(W, X)$, the k -th component of the output vector of $G(W, X)$ is interpreted as the “cost” of classifying X into category k .
- Best output: $\check{Y} = \min_{Y \in \{Y\}} E(W, Y, X) = \min_k G_k(W, X)$.

Examples of EBM Classifier: Perceptron



- Y is a discrete variable, $\{Y\} = \{-1, +1\}$.
- Energy: $E(W, Y, X) = -Y.W'X$.
- Best output: $\check{Y} = \text{sign}(W'X)$, where $\text{sign}(R) = +1$ iff $R > 0$ and -1 otherwise.

Linear Machines



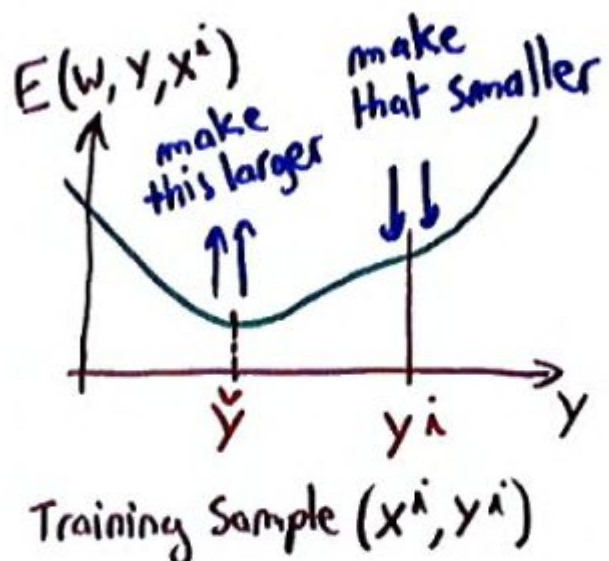
- The learning algorithms we have seen so far (perceptron, linear regression) are of that form, with the assumption that $G(W, X)$ only depends on the dot product of W and X .
- In other words, The E function of 2-class linear classifiers can be written as:

$$E(Y, X, W) = D(Y, f(W'X))$$

where $W'X$ is the dot product of vectors W and X , and f is a monotonically increasing scalar function.

- in the following, we assume $Y = -1$ for class 1, and $Y = +1$ for class 2.

Training Energy-Based Models



- To train an EBM, we minimize a **loss function**, which is an average over training samples of a **per-sample loss function** $L(W, Y^i, X^i)$:

$$\mathcal{L}(W, \mathcal{S}) = \frac{1}{P} \sum_{i=1}^P L(W, Y^i, X^i)$$

- The loss function must be designed so that minimizing it with respect to W will make the machine approach the desired behavior.

To ensure this, we pick loss functions that, for a given training input X^i , will drive the energies $E(W, Y^i, X^i)$ associated with the desired output Y^i to be lower than the energies associated with all other (undesired) outputs values $E(W, Y, X^i)$ for all $Y \neq Y^i, Y \in \{Y\}$.

Form of the Loss Function

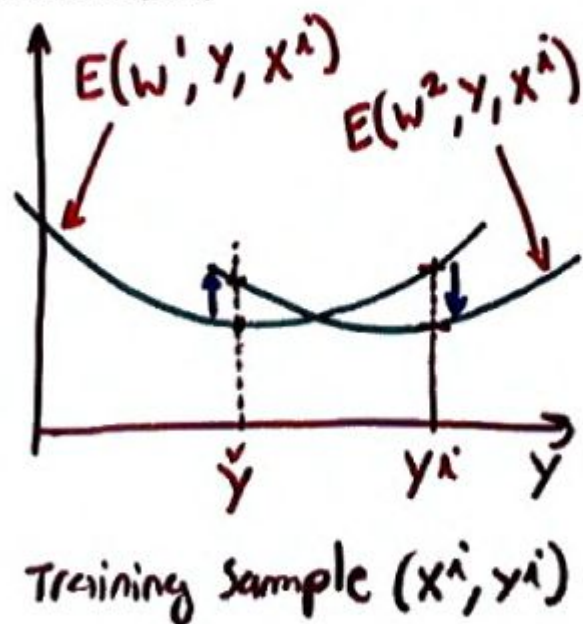
- We assume that the per-sample loss function $L(W, Y^i, X^i)$ has a lower bound over W for all Y^i, X^i .
- We assume that L depends on X^i only indirectly through the set of energies $\{E(W, Y, X^i), Y \in \{Y\}\}$.
- For example, if $\{Y\}$ is the set of integers between 0 and $k - 1$ (as would be the case for a classifier with k categories), the per-sample loss for sample (X^i, Y^i) should be of the form:

$$L(W, Y^i, X^i) = L(Y^i, E(W, 0, X^i), E(W, 1, X^i), \dots, E(W, k - 1, X^i))$$

- With this assumption, we separate the choice of the loss function from the details of the internal structure of the machine, and limit the discussion to how minimizing the loss function affects the energies.

Examples of Loss: Energy Loss

Energy Loss, the simplest of all losses: $L_{\text{energy}}(W, Y^i, X^i) = E(W, Y^i, X^i)$. This loss only works if $E(W, Y, X^i)$ has a special form which guarantees that making $E(W, Y^i, X^i)$ lower will automatically make $E(W, Y, X^i)$ for $Y \neq Y^i$ larger than the minimum.

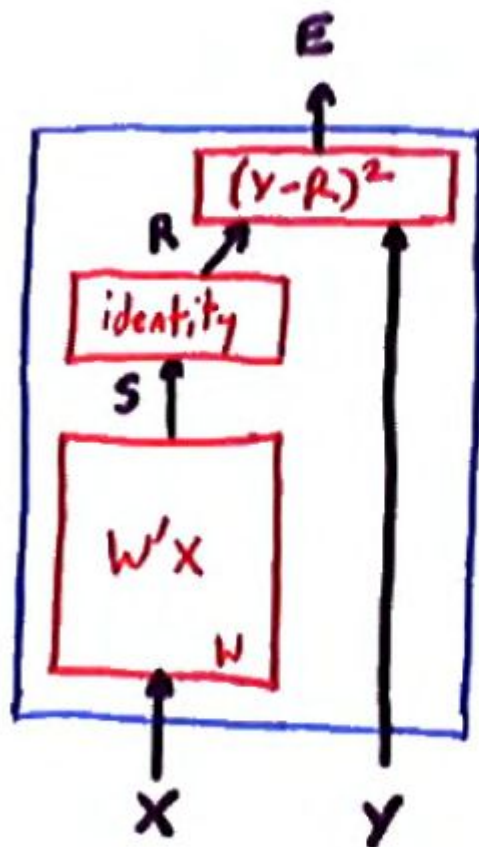


Example: if $E(W, Y, X)$ is quadratic in Y , as is the case for regression with squared error: $E(W, Y, X) = \|Y - G(W, X)\|^2$,
Let $W(1)$ is the parameter before a learning update, and $W(2)$ the parameter after the learning update, and let $\check{Y} = \min_Y E(W(1), Y, X)$. Then,

$$E(W(2), Y^i, X^i) - E(W(2), \check{Y}, X^i) < E(W(1), Y^i, X^i) - E(W(1), \check{Y}, X^i)$$

Linear Regression

Linear regression uses the Energy loss

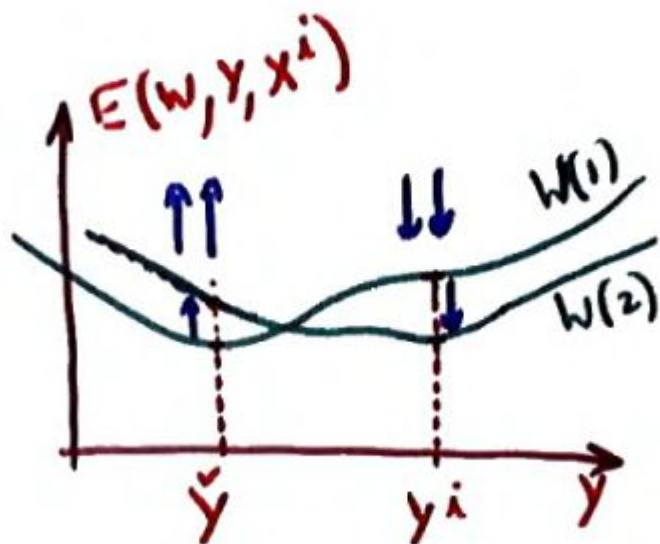


- $R = W'X$
- $E(W, Y, X) = D(Y, R) = \frac{1}{2} \|Y - R\|^2$
- $L(W, Y^i, X^i) = D(Y^i, W'X^i)$
- $\frac{\partial L}{\partial W} = \frac{\partial D(Y^i, R)}{\partial R} \frac{\partial R}{\partial W}$
- $\frac{\partial L}{\partial W} = \frac{\partial D(Y^i, R)}{\partial R} \frac{\partial (W'X^i)}{\partial W} = (R - Y^i)X^i$
- descent: $W \leftarrow W + \eta(Y^i - R)X^i$

Examples of Loss: Perceptron Loss

Perceptron Loss:

$$L_{\text{perceptron}}(W, Y^i, X^i) = E(W, Y^i, X^i) - \min_{Y \in \{Y\}} E(W, Y, X^i)$$

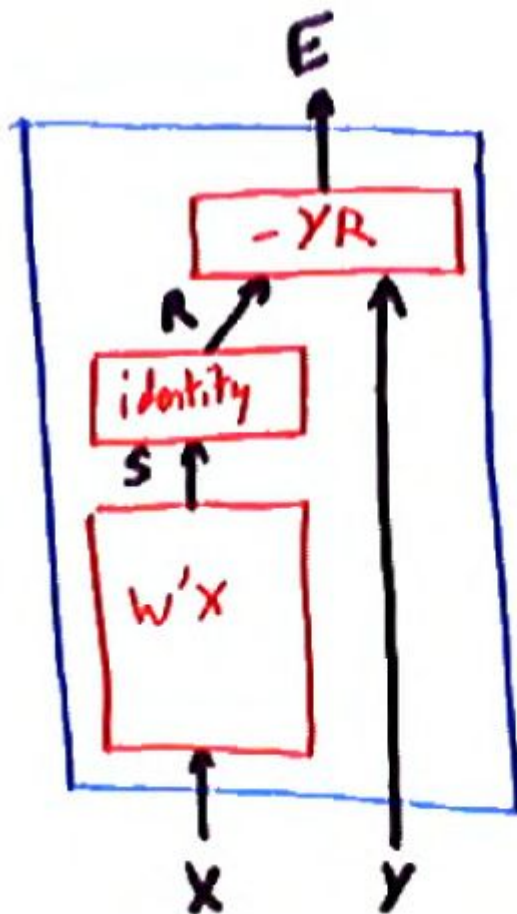


Adjust W so that $E(W, Y^i, X^i)$ gets smaller, while $\check{Y} = \min_{Y \in \{Y\}} E(W, Y, X^i)$ gets bigger (or more precisely, so that the difference decreases). This algorithm makes no update whenever the energy of the desired Y is lower than all the others.

Perceptron

$$L_{\text{perceptron}}(W, Y^i, X^i) = E(W, Y^i, X^i) - \min_{Y \in \{Y\}} E(W, Y, X^i)$$

$$\{Y\} = \{-1, +1\}.$$



- $R = W'X$
- $E(Y, X, W) = D(Y, R) = -YR$
- $Y \in \{-1, +1\}$, hence $\min_Y -YR = -\text{sign}(R)R$ where $\text{sign}(R) = 1$ iff $R > 0$, and -1 otherwise.
- $L(W, Y^i, X^i) = -(Y^i - \text{sign}(R))R$
- $\frac{\partial L}{\partial W} = \frac{\partial -(Y^i - \text{sign}(R))R}{\partial R} \frac{\partial R}{\partial W}$
- $\frac{\partial L}{\partial W} = -(Y^i - \text{sign}(W'X^i))X^i$
- descent: $W \leftarrow W + \eta(Y^i - \text{sign}(W'X^i))X^i$

Examples of Loss: Log-Likelihood Loss

Log-Likelihood Loss:

$$L_{ll}(W, Y^i, X^i) = E(W, Y^i, X^i) + \frac{1}{\beta} \log \left(\sum_{Y \in \{Y\}} \exp(-\beta E(W, Y, X^i)) \right)$$

where β is a positive constant.

- The function $\mathcal{F}_\beta(\{Y\}) = \frac{1}{\beta} \log \left(\sum_{Y \in \{Y\}} \exp(-\beta E(W, Y, X^i)) \right)$ is called the **free energy** of the ensemble $\{Y\}$ for temperature $1/\beta$.

- We define $\mathcal{Z}_\beta(\{Y\}) = \sum_{Y \in \{Y\}} \exp(-\beta E(W, Y, X^i))$ as the **partition function** of ensemble $\{Y\}$.

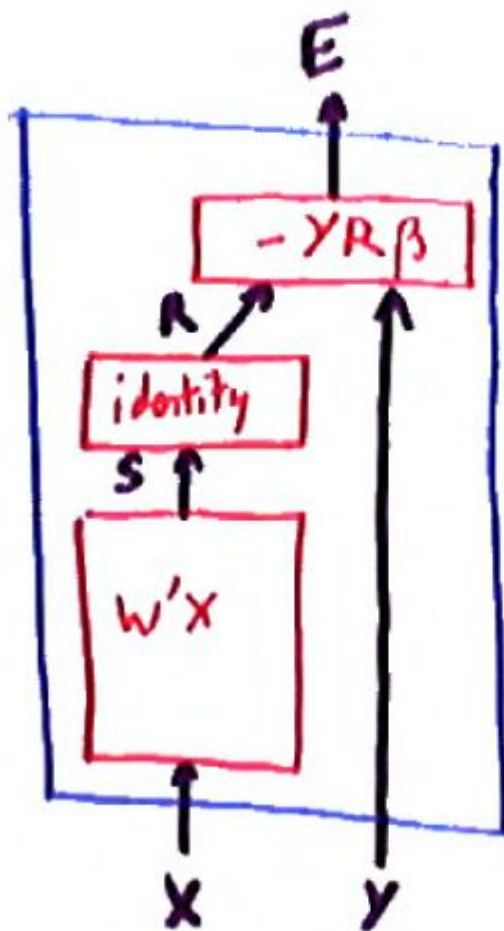
- Interesting property # 1: $\mathcal{F}_\beta(\{Y\}) = \frac{1}{\beta} \log \mathcal{Z}_\beta(\{Y\})$

- Interesting property # 2: $\lim_{\beta \rightarrow \infty} \mathcal{F}_\beta(\{Y\}) = \min_{Y \in \{Y\}} E(W, Y, X^i)$

For very large β , the log-likelihood loss reduces to the Perceptron loss.

Logistic Regression (a.k.a MaxEnt)

$$L_{ll}(W) = E(Y^i, X^i, W) + \log \left(\sum_{Y \in \{Y\}} \exp(-E(W, Y, X^i)) \right)$$



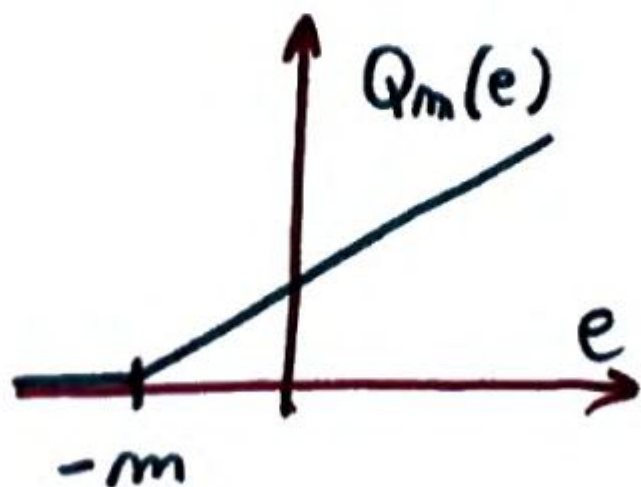
- $R = \frac{1}{2} W' X$
- $E(Y, X, W) = D(Y, R) = -\frac{1}{2} Y R = -\frac{1}{2} Y W' X$
- $L(W) = \log(1 + \exp(-Y^i W' X^i))$
- $\frac{\partial L}{\partial W} = \frac{\partial D(Y^i, R)}{\partial R} \frac{\partial S}{\partial W}$
- $\frac{\partial L}{\partial W} = - \left(\frac{Y^i + 1}{2} - \frac{1}{1 + \exp(-W' X^i)} \right) X^i$
- descent: $W \leftarrow W + \eta \left(\frac{Y^i + 1}{2} - \frac{1}{1 + \exp(-W' X^i)} \right) X^i$

Examples of Loss: Margin Loss

Margin Loss: for discrete output set $\{Y\}$:

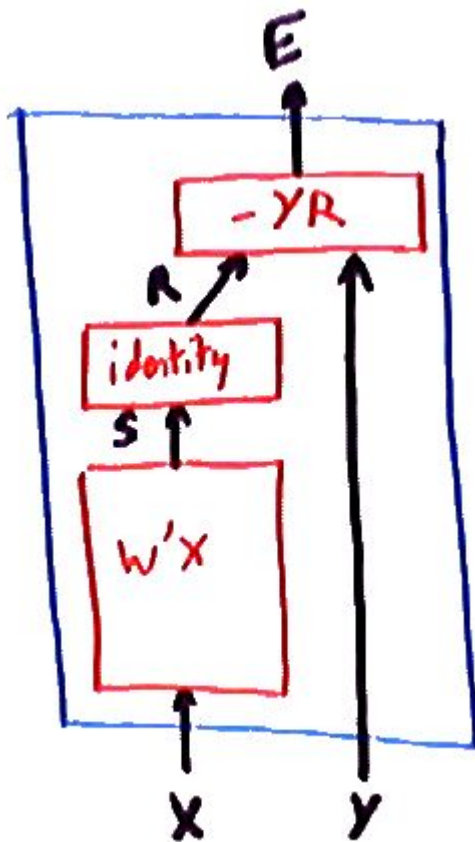
$$L_{\text{margin}}(W, Y^i, X^i) = Q_m \left(E(W, Y^i, X^i) - \min_{Y \in \{Y\}, Y \neq Y^i} E(W, Y, X^i) \right)$$

where $Q_m(e)$ is any function that is monotonically increasing for $e > -m$, where m is a constant called the **margin**.

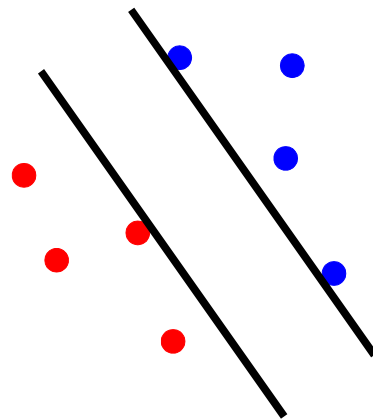


Adjust W so that $E(W, Y^i, X^i)$ gets smaller, while all $E(W, Y, X^i)$ for which $E(W, Y, X^i) - E(W, Y^i, X^i) < m$ get bigger. This guarantees that the energy of the desired Y will be smaller than all other energies by at least m .

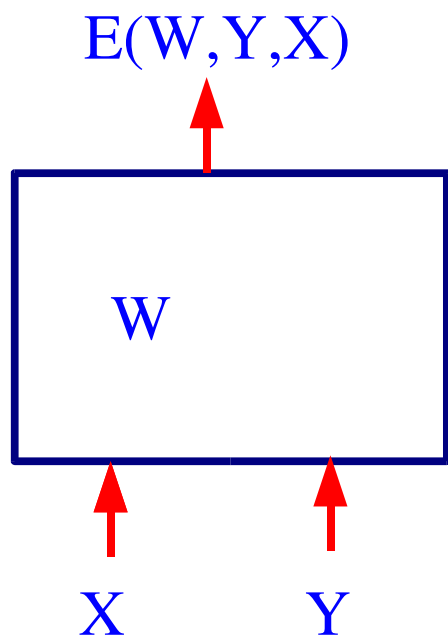
Linear Model + Margin Loss + Regularization = SVM



- Minimize the hinge loss: make the energy of all the “good” answers smaller than the energy of any “bad” answer by at least m (the margin).
- Minimize the Regularization term: Make W as short as possible.
- This is equivalent to keeping $\|W\|$ constant, while maximizing m .

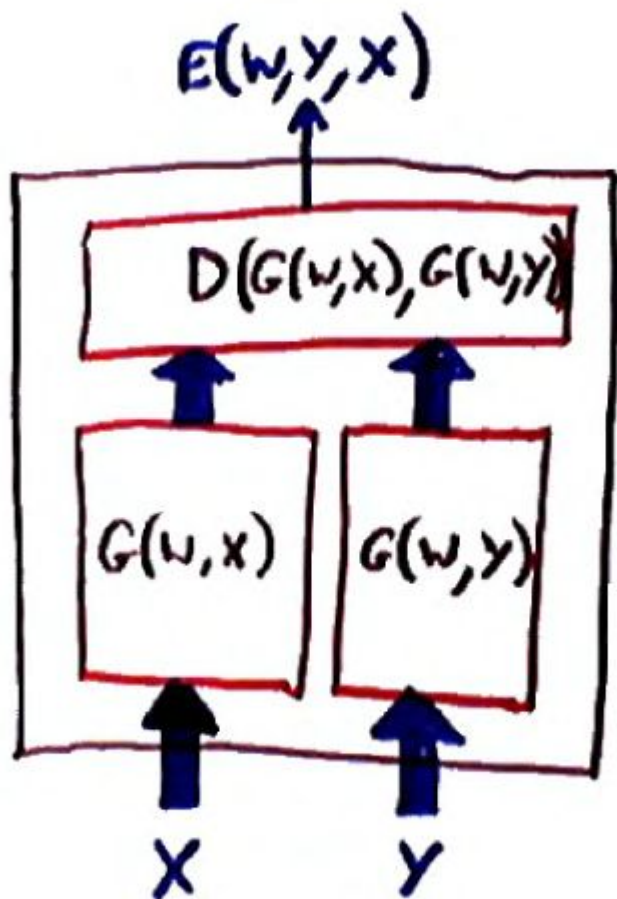


Architecture



- We can put anything we want in the box.
- The energy can be a **very complicated non-linear function of X, Y, and W** (e.g. **A neural net, a graphical model, an HMM, Markov Random Field,....**).
- **The internal structure of the box is called the architecture of the model.**

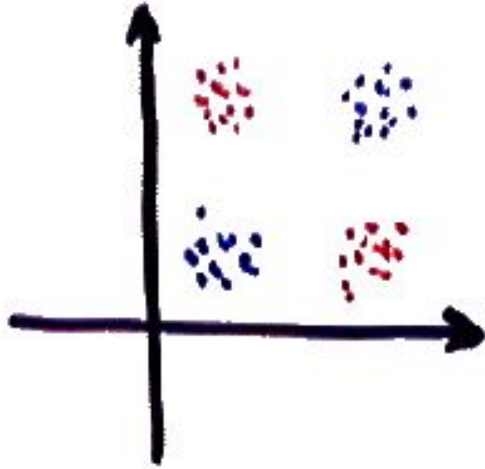
Examples of EBM: Matcher



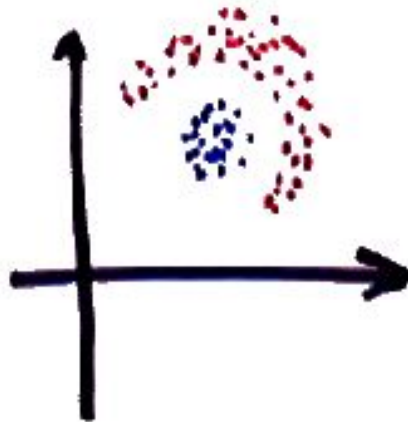
- X and Y are vectors of the same dimension.
- Energy:
 $E(W, Y, X) = D(G(W, Y), G(W, X))$ where $D(.,.)$ is a distance or dissimilarity measure.
- Best output: $\check{Y} = \min_Y E(W, Y, X) = G^{(-1)}(G(W, X))$.

Finding the Y that minimizes the energy may be non-trivial

Limitations of Linear Machines

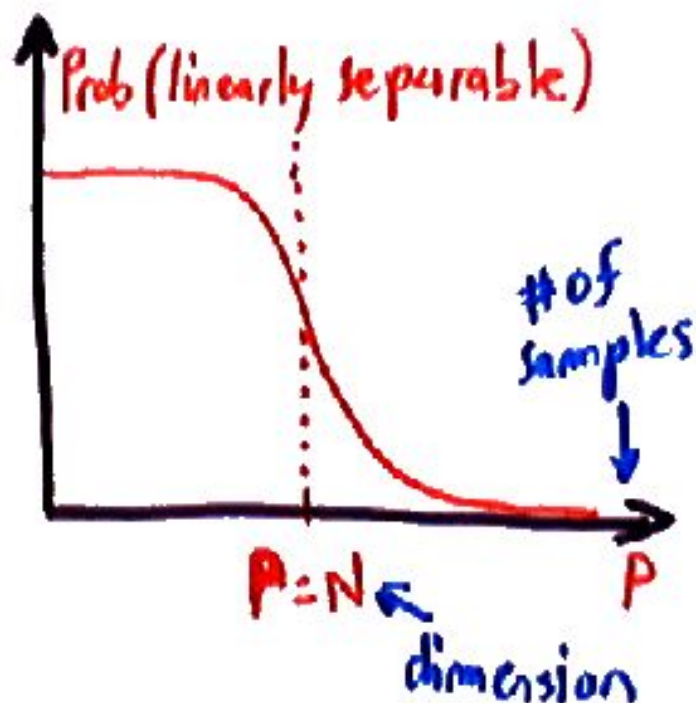


The *Linearly separable* dichotomies are the partitions that are realizable by a linear classifier (the boundary between the classes is a hyperplane).



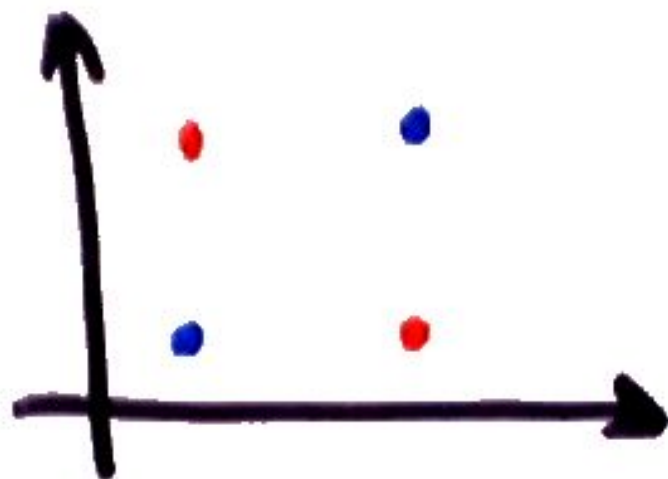
Number of Linearly Separable Dichotomies

The probability that P samples of dimension N are linearly separable goes to zero very quickly as P grows larger than N (Cover's theorem, 1966).

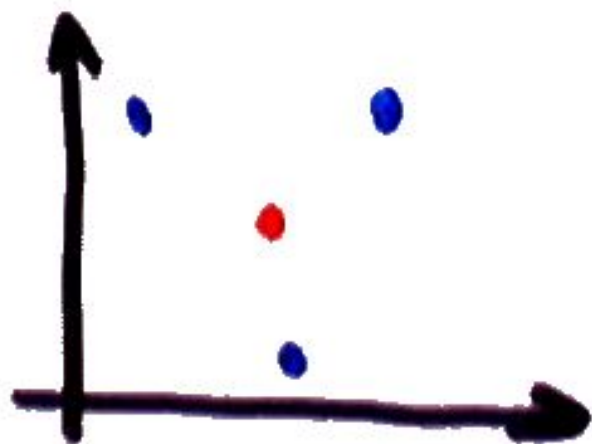


- Problem: there are 2^P possible dichotomies of P points.
- Only about N are linearly separable.
- If P is larger than N , the probability that a random dichotomy is linearly separable is very, very small.

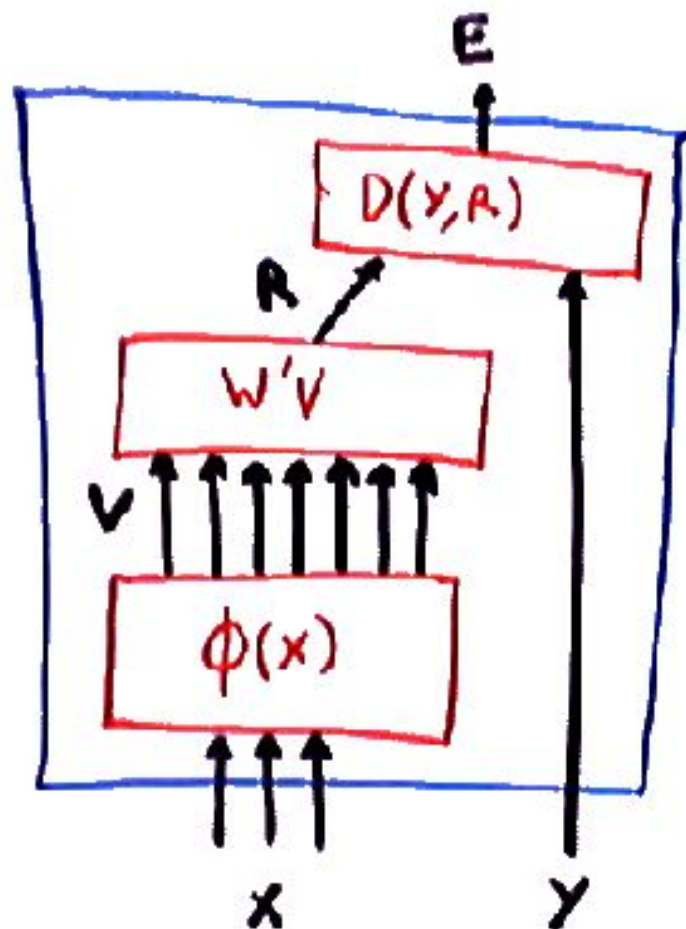
Example of Non-Linearly Separable Dichotomies



- Some seemingly simple dichotomies are not linearly separable
- **Question:** How do we make a given problem linearly separable?

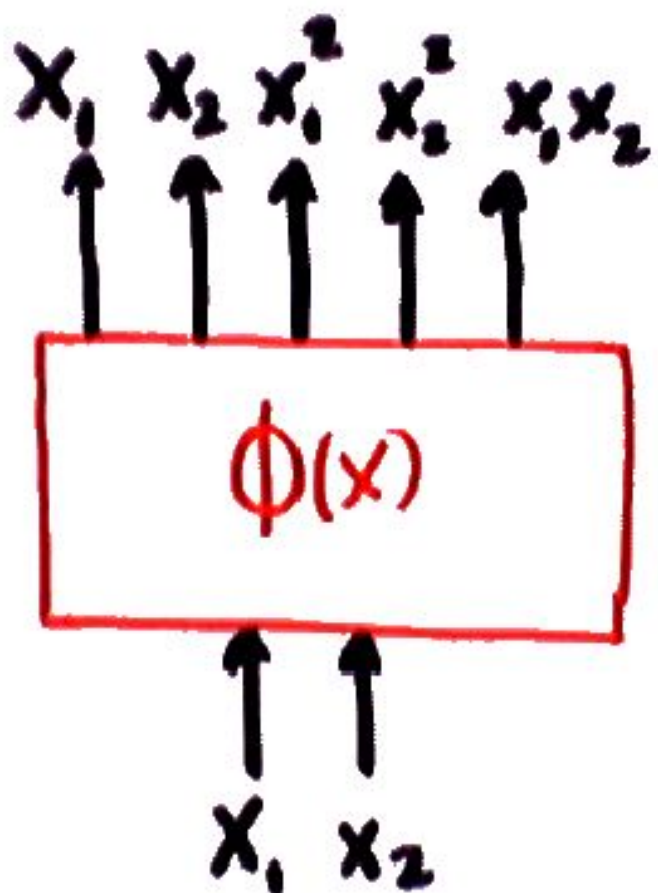


Making N Larger: Preprocessing



- **Answer 1:** we make N larger by augmenting the input variables with new “features”.
- we map/project X from its original N -dimensional space into a higher dimensional space where things are more likely to be linearly separable, using a vector function $\Phi(X)$.
- $E(Y, X, W) = D(Y, R)$
- $R = f(W'V)$
- $V = \Phi(X)$

Adding Cross-Product Terms



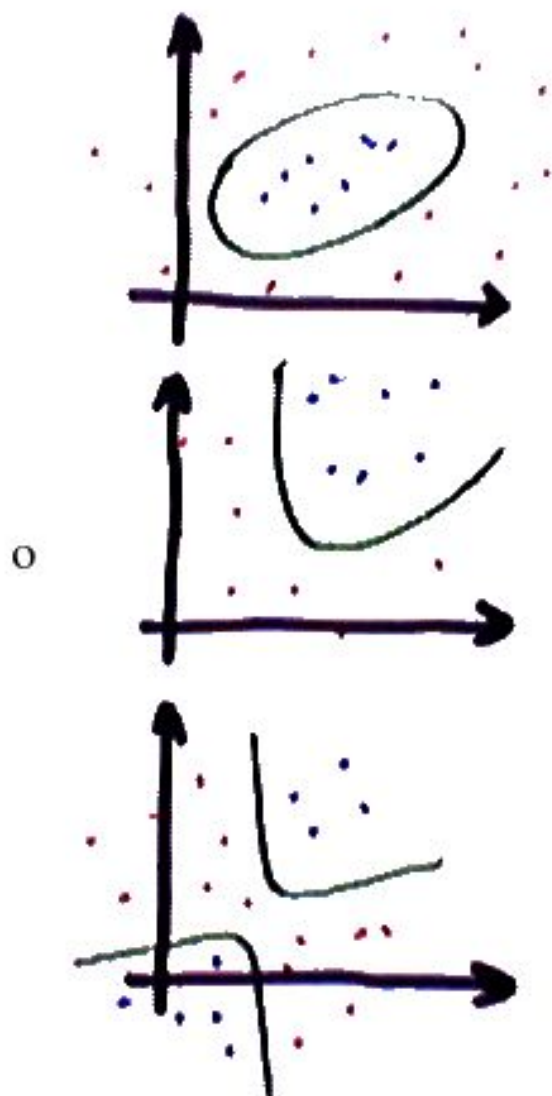
- Polynomial Expansion.
- If our original input variables are $(1, x_1, x_2)$, we construct a new *feature vector* with the following components:

$$\Phi(1, x_1, x_2) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$$

i.e. we add all the cross-products of the original variables.

- we map/project X from its original N -dimensional space into a higher dimensional space with $N(N+1)/2$ dimensions.

Polynomial Mapping



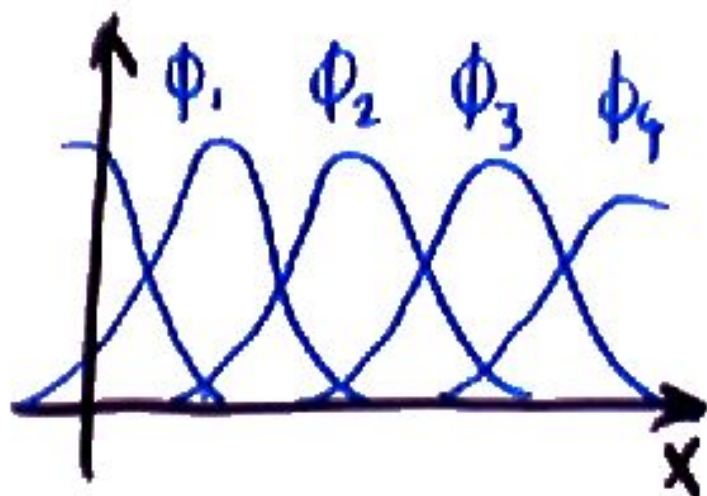
- Many new functions are now separable with the new architecture.
- With cross-product features, the family of class boundaries in the original space is the conic sections (ellipse, parabola, hyperbola).
- to each possible boundary in the original space corresponds a linear boundary in the transformed space.
- Because this is essentially a linear classifier with a preprocessing, we can use standard linear learning algorithms (perceptron, linear regression, logistic regression...).

Problems with Polynomial Mapping

- We can generalize this idea to higher degree polynomials, adding cross-product terms with 3, 4 or more variables.
- Unfortunately, the number of terms is the number of combinations d choose N , which grows like N^d , where d is the degree, and N the number of original variables.
- In particular, the number of free parameters that must be learned is also of order N^d .
- This is impractical for large N and for $d > 2$.
- Example: handwritten digit recognition (16x16 pixel images). Number of variables: 256. Degree 2: 32,896 variables. Degree 3: 2,796,160. Degree 4: 247,460,160.....

Next Idea: Tile the Space

place a number of equally-spaced “bumps” that cover the entire input space.



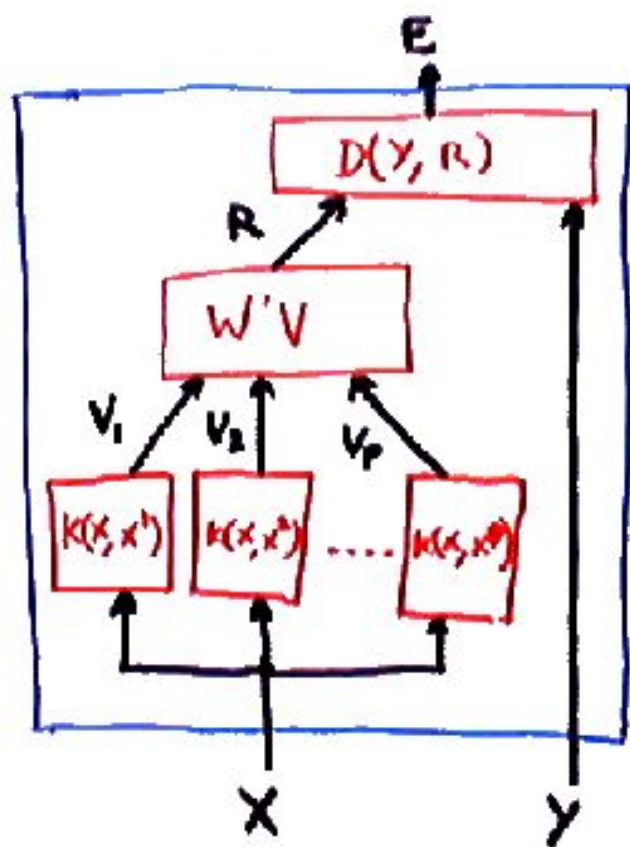
- For classification, the bumps can be Gaussians
- For regression, the basis functions can be wavelets, sine/cosine, splines (pieces of polynomials)....
- **problem:** this does not work with more than a few dimensions.
- The number of bumps necessary to cover an N dimensional space grows exponentially with N .

Sample-Centered Basis Functions (Kernels)

Place the center of a basis function around each training sample. That way, we only spend resources on regions of the space where we actually have training samples.

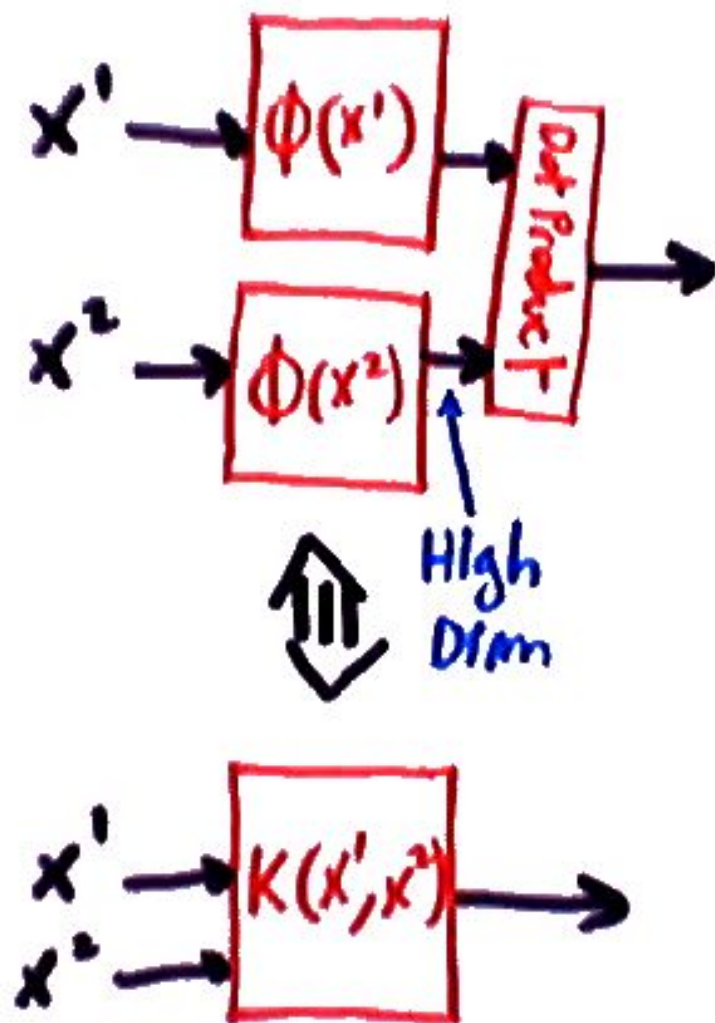
- Discriminant function:

$$f(X, W) = \sum_{k=1}^{k=P} W_k K(X, X^k)$$



- $K(X, X')$ often takes the form of a *radial basis function*:
 $K(X, X') = \exp(b\|X - X'\|^2)$ or a polynomial $K(X, X') = (X \cdot X' + 1)^m$
- This is a very common architecture, which can be used with a number of energy functions.
- In particular, this is the architecture of the so-called **Support Vector Machine** (SVM), but the energy function of the SVM is a bit special. We will study it later in the course.

The Kernel Trick



- If the kernel function $K(X, X')$ verifies the *Mercer conditions*, then there exist a mapping Φ , such that $\Phi(X) \cdot \Phi(X') = K(X, X')$.
- The Mercer conditions are that K must be symmetric, and must be positive definite (i.e $K(X, X)$ must be positive for all X).
- In other words, if we want to map our X into a high-dimensional space (so as to make them linearly separable), and all we have to do in that space is compute dot products, we can take a shortcut and simply compute $K(X^1, X^2)$ without going through the high-dimensional space.
- This is called the “**kernel trick**”. It is used in many so-called Kernel-based methods, including Support Vector Machines.

Examples of Kernels

- **Quadratic kernel:** $\Phi(X) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2)$ then

$$K(X, X') = \Phi(X) \cdot \Phi(X') = (X \cdot X' + 1)^2$$

- **Polynomial kernel:** this generalizes to any degree d . The kernel that corresponds to $\Phi(X)$ being a polynomial of degree d is

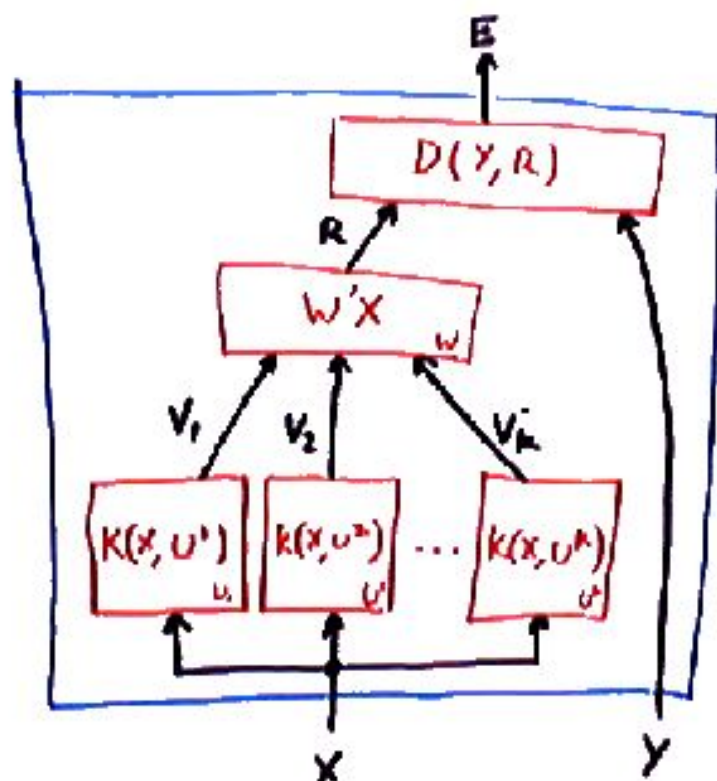
$$K(X, X') = \Phi(X) \cdot \Phi(X') = (X \cdot X' + 1)^d.$$

- **Gaussian Kernel:**

$$K(X, X') = \exp(-b\|X - X'\|^2)$$

This kernel, sometimes called the Gaussian Radial Basis Function, is very commonly used.

Sparse Basis Functions

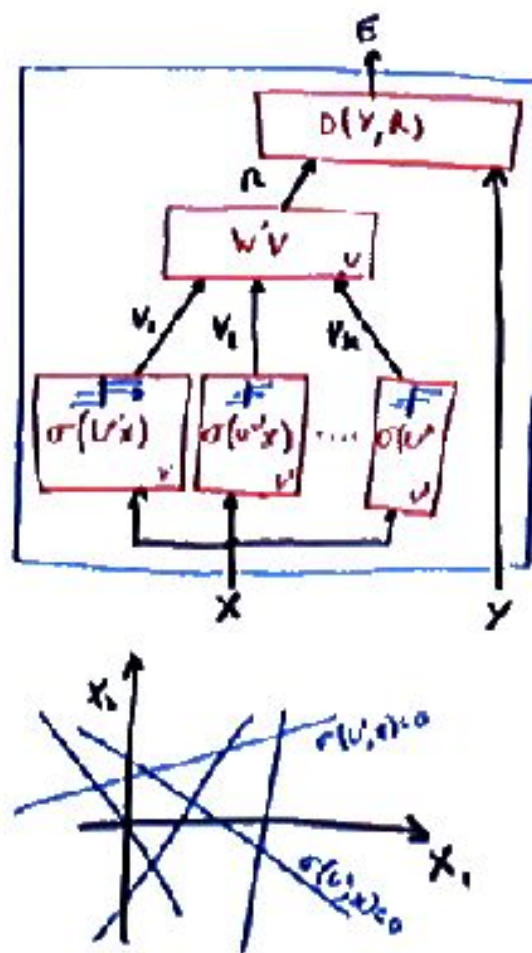


- Place the center of a basis function around areas containing training samples.
- Idea 1: use an unsupervised clustering algorithm (such as K-means or mixture of Gaussians) to place the centers of the basis functions in areas of high sample density.
- Idea 2: adjust the basis function centers through gradient descent in the loss function.

The discriminant function F is:

$$F(X, W, U^1, \dots, U^K) = \sum_{k=1}^{k=K} W_k K(X, U^k)$$

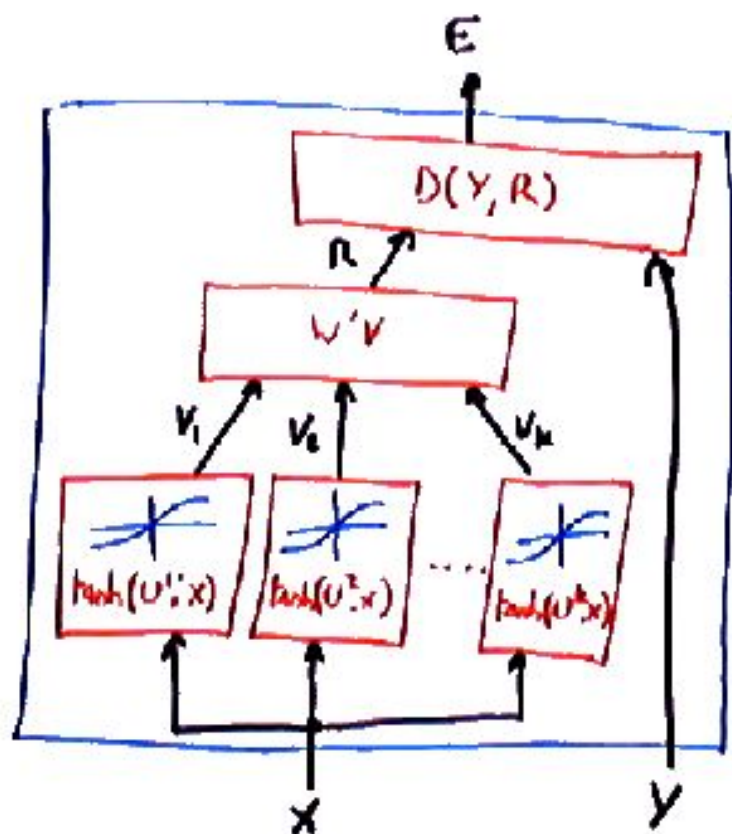
Other Idea: Random Directions



- Partition the space in lots of little domains by randomly placing lots of hyperplanes.
- Use many variables of the type $q(W^k X)$, where q is the threshold function (or some other squashing function) and W_k is a randomly picked vector.
- This is the original Perceptron.
- Without the non-linearity, the whole system would be linear (product of linear operations), and therefore would be no more powerful than a linear classifier.
- **problem:** a bit of a wishful thinking, but it works occasionally.

Neural Net with a Single Hidden Layer

A particularly interesting type of basis function is the sigmoid unit: $V_k = \tanh(U'^k X)$



- a network using these basis functions, whose output is $R = \sum_{k=1}^{K} W_k V_k$ is called a *single hidden-layer neural network*.
- Similarly to the RBF network, we can compute the gradient of the loss function with respect to the U^k :

$$\begin{aligned} \frac{\partial L(W)}{\partial U^j} &= \frac{\partial L(W)}{\partial R} W_j \frac{\partial \tanh(U^j X)}{\partial U^j} \\ &= \frac{\partial L(W)}{\partial R} W_j \tanh'(U^j X) X' \end{aligned}$$

Any well-behaved function can be approximated as close as we wish by such networks (but K might be very large).

