

Introduction to Neural Networks

Supervised learning: neural networks, statistics, & machine learning

Initialize:

```
In[699]:= Off[General::spell1];  
  
In[700]:= SetOptions[ListPlot, ImageSize -> Small];  
          SetOptions[DensityPlot, ImageSize -> Small];  
          SetOptions[ContourPlot, ImageSize -> Small];
```

Overview

Previously...

Introduced the notions of loss and risk for inference.

Rationale for Bayesian from the point of view of supervised learning that minimizes empirical risk.

We did this by starting from a goal to learn a decision rule (or estimator) to minimize average loss over a set of training pairs (the empirical risk), and then showed that this is equivalent (in the limit) to making a Bayes optimal decision. This provides a tie between empirical methods such as support vector machines discussed in this lecture, and Bayesian decision theory.

Today

Set neural network supervised learning in the context of various statistical/machine learning methods.

Rationale: The relationship of brain to behavior is complicated. Good to understand bottom-up, from neurons to behavior. But also good to understand top-down, from behavior to quantitative models with as few free parameters as possible. But with a view to plausible neural network functioning.

In other words, always a good idea to try the simplest model first.

Preview of “big data” and scientific python

Nearest neighbor demonstration: “random Xor”

Just because a problem isn't linearly separable, doesn't mean the algorithm has to be complicated. The nearest neighbor classifier is one of the simplest algorithms for supervised learning. If one has a prior rationale for a representation, and measure of distance (or similarity) that reflects how similar patterns cluster based on category, the nearest neighbor algorithm is a good thing to try.

k nearest neighbors or k-NN for short. Has no problem with Xor.

Define \mathcal{D} to be a mixture distribution of bivariate normals

```
In[703]:=  $\mathcal{D}$  = MixtureDistribution[{1, 1, 1, 1},
  {MultinormalDistribution[{0, 0.}, {{1., 0}, {0, 1.}}],
  MultinormalDistribution[{0, 10.}, {{1., 0}, {0, 1.}}],
  MultinormalDistribution[{10, 0.}, {{1., 0}, {0, 1.}}],
  MultinormalDistribution[{10, 10.}, {{1., 0}, {0, 1.}}]}];
```

Here's a sample:

```
In[704]:= RandomVariate[ $\mathcal{D}$ ]
Out[704]:= {10.5241, 8.90777}
```

We can generate $n = 100$ samples:

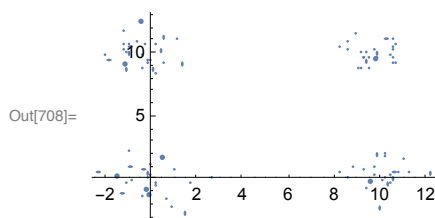
```
In[705]:= n = 100;
In[706]:= somedata = Table[RandomVariate[ $\mathcal{D}$ ], {n}];
```

And then artificially assign labels, either 0 or 1 to each of the 100 points.

► 1. Exercise: change below for Xor, Or...And...

```
In[707]:= labels = Which[EuclideanDistance[#, {0, 0}] < 5, 0,
  EuclideanDistance[#, {0, 10}] < 5, 1, EuclideanDistance[#, {10, 0}] < 5,
  1, EuclideanDistance[#, {10, 10}] < 5, 0] & /@ somedata;
```

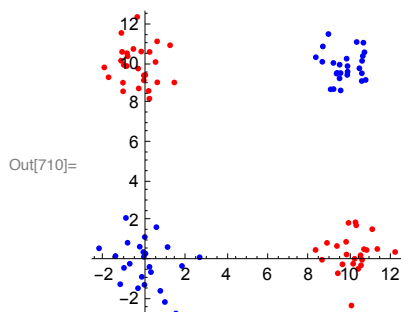
```
In[708]:= ListPlot[somedata]
```



We can color the plots with the labels, red for 1's and blue for 0's.

```
In[709]:= colorf = Blend[{{0, Blue}, {1, Red}}, #] &
pl = Graphics[MapThread[{colorf[#1], Point[#2]} &, {labels, somedata}],
  Axes → True, AspectRatio → 1]
```

```
Out[709]= Blend[{{0, Blue}, {1, Red}}, #1] &
```



```
In[711]:= Dimensions[somedata]
```

```
Out[711]:= {100, 2}
```

Estimate nearest neighbor boundaries

OK, so we have our artificial data and a pretty plot. Let's classify.

Run through lots of points inside $\{0,10\} \times \{0,10\}$. For each point find the k nearest neighbors, and count how many have labels of 1 (i.e. how many dots are "red"). Plot up the count proportions--i.e. plot the proportion of 0s or 1s falling within a circle containing $k=5$ neighbors, for each (x,y) point. ("circle" is implicit in the default distance function used by `Nearest[]`, which has `EuclideanDistance` as the default distance measure).

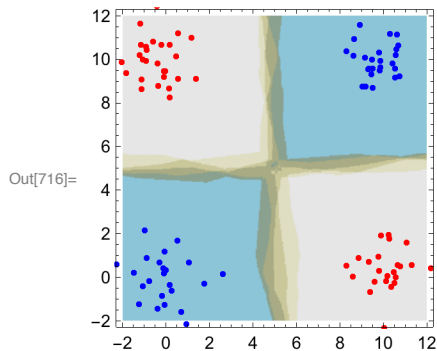
```
In[712]:= trainingdata = Thread[somedata → labels];
```

```
In[713]:= Clear[nf2];
```

```
nf2[x_, k_] := Total[Nearest[trainingdata, x, k]] / k;
```

```
In[715]:= k = 5;
```

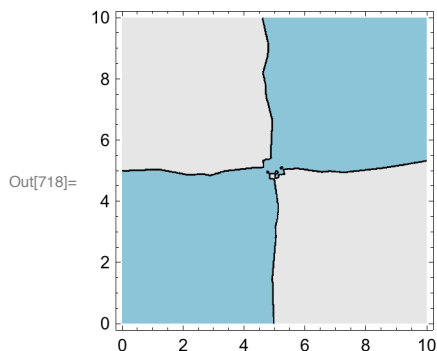
```
Show[ {DensityPlot[nf2[{x, y}, k], {x, -2, 12}, {y, -2, 12},
  PlotPoints -> n, ColorFunction -> "LightTerrain", Exclusions -> None], pl}]
```



Make the decision: if the proportion is more than $1/2$ decide 1, otherwise 0:

```
In[717]:= k = 5;
```

```
ContourPlot[If[nf2[{x, y}, k] > 1/2, 1, 0], {x, 0, 10}, {y, 0, 10},
  PlotPoints -> 60, ColorFunction -> "LightTerrain", Exclusions -> None]
```



$k=1$ and Voronoi diagrams

How to pick K ? Note that small values of k produce high variance predictors--variations in future samples produce big differences in the decision boundary. If k includes all the training data, the prediction is

the most frequently occurring label--i.e. behaves like it is basing its decision on the prior. A solution is to use cross-validation, where one divides the training pairs into a training set and a test or validation set. And then compare the results of training with various values of k on the training set with respect to the ability to classify members of the test or validation set.

Nearest neighbor regression

"If something is similar to something else in one respect, it is likely to be similar in another respect." Can be OK for smooth interpolation, but can have problem for extrapolation.

Return to linear separability

Discriminant functions

Before looking at traditional (Fisher) and more modern discriminants (linear and non-linear SVM), let's build our geometric intuitions of what a simple threshold logic or perceptron unit does by viewing it from a more formal point of view. Perceptron learning is an example of nonparametric statistical learning, because it doesn't require knowledge of the underlying probability distributions generating the data (such distributions are characterized by a relatively small number of "parameters", such as the mean and variance of a Gaussian distribution). Of course, how well it does will depend on the generative structure of the data. Much of the material below is covered in Duda and Hart (1978).

Linear discriminant functions: Two category case

Let \mathbf{x} be a feature vector. A discriminant function, $g(\mathbf{x})$ divides input space into two category regions depending on whether $g(\mathbf{x}) > 0$ or $g(\mathbf{x}) < 0$. (We've switched notation, $\mathbf{x}=\mathbf{f}$). This is the hard threshold case, i.e. a step function.

This linear case then corresponds to the simple perceptron unit:

$$g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$$

where \mathbf{w} is the weight vector and w_0 is the (scalar) threshold (which we've called bias, although this "bias" has nothing to do with statistical "bias"), with the decision: $g(\mathbf{x}) > 0$ or $g(\mathbf{x}) < 0$, corresponding to a step-function non-linearity.

Discriminant functions can be generalized, for example to quadratic decision surfaces:

$$g(\mathbf{x}) = w_0 + \sum_{i=1} w_i x_i + \sum_{i=1} \sum_{j=1} w_{ij} x_i x_j$$

where $\mathbf{x} = \{x_1, x_2, x_3, \dots\}$. $g(\mathbf{x})=0$ defines a decision surface. For now we restrict ourselves to the linear case where the decision surface is a hyperplane.

Suppose \mathbf{x}_1 and \mathbf{x}_2 are vectors, with endpoints sitting on the hyperplane (or in the simple case of the red line below), then their difference is a vector lying in the hyperplane

$$\begin{aligned} g(\mathbf{x}_1) &= g(\mathbf{x}_2) = \mathbf{w} \cdot \mathbf{x}_1 + w_0 = \mathbf{w} \cdot \mathbf{x}_2 + w_0 \\ \mathbf{w} \cdot (\mathbf{x}_1 - \mathbf{x}_2) &= 0 \end{aligned}$$

so because the dot product is zero, the weight vector \mathbf{w} is normal to any vector lying in the hyperplane.

Thus \mathbf{w} determines how the plane is oriented. The normal vector \mathbf{w} points into the region for which $g(\mathbf{x}) > 0$, and $-\mathbf{w}$ points into the region for which $g(\mathbf{x}) < 0$.

$$\mathbf{w} \cdot \mathbf{x} + w_0 = 0, \text{ so } \mathbf{w} \cdot \mathbf{x} = -w_0.$$

If we project \mathbf{x} onto the normalized weight vector $\mathbf{w}/|\mathbf{w}|$, we have the normal distance of the hyperplane from the origin equal to:

$$\mathbf{w} \cdot \mathbf{x} / |\mathbf{w}| = -w_0 / |\mathbf{w}|$$

Thus, the threshold determines the position of the hyperplane.

One can also show that the normal distance of an arbitrary vector \mathbf{x} to the hyperplane is given by:

$$g(\mathbf{x}) / |\mathbf{w}|$$

► 2. Prove $g(\mathbf{x})/|\mathbf{w}|$ is the distance of a vector \mathbf{x} to the decision hyperplane

Let \mathbf{x}_p be the vector to the point of projection on the discriminant. Then $(\mathbf{x} - \mathbf{x}_p)$ is the vector from \mathbf{x} to \mathbf{x}_p . If we project $(\mathbf{x} - \mathbf{x}_p)$ on to $\mathbf{w}/|\mathbf{w}|$ (which is the unit vector perpendicular to the decision plane), then $(\mathbf{x} - \mathbf{x}_p) \cdot \mathbf{w}/|\mathbf{w}|$ is the distance we want.

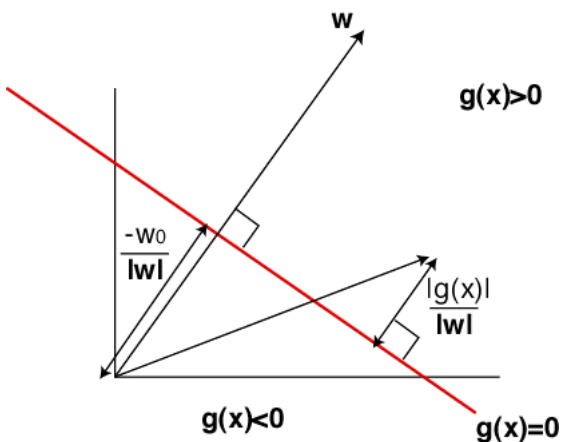
$$\mathbf{x} \cdot \mathbf{w} - \mathbf{x}_p \cdot \mathbf{w} = \mathbf{x} \cdot \mathbf{w} - g(\mathbf{x}_p) + w_0 = \mathbf{x} \cdot \mathbf{w} + w_0 = g(\mathbf{x})$$

So the distance is $g(\mathbf{x})/|\mathbf{w}|$.

To summarize:

- 1) discriminant function divides the input space by a hyperplane decision surface;
- 2) The orientation of the surface is determined by the weight vector \mathbf{w} ;
- 3) the location is determined by the threshold w_0 ;
- 4) the discriminant function gives a measure of how far an input vector is from the hyperplane.

The figure summarizes the basic properties of the linear discriminant.

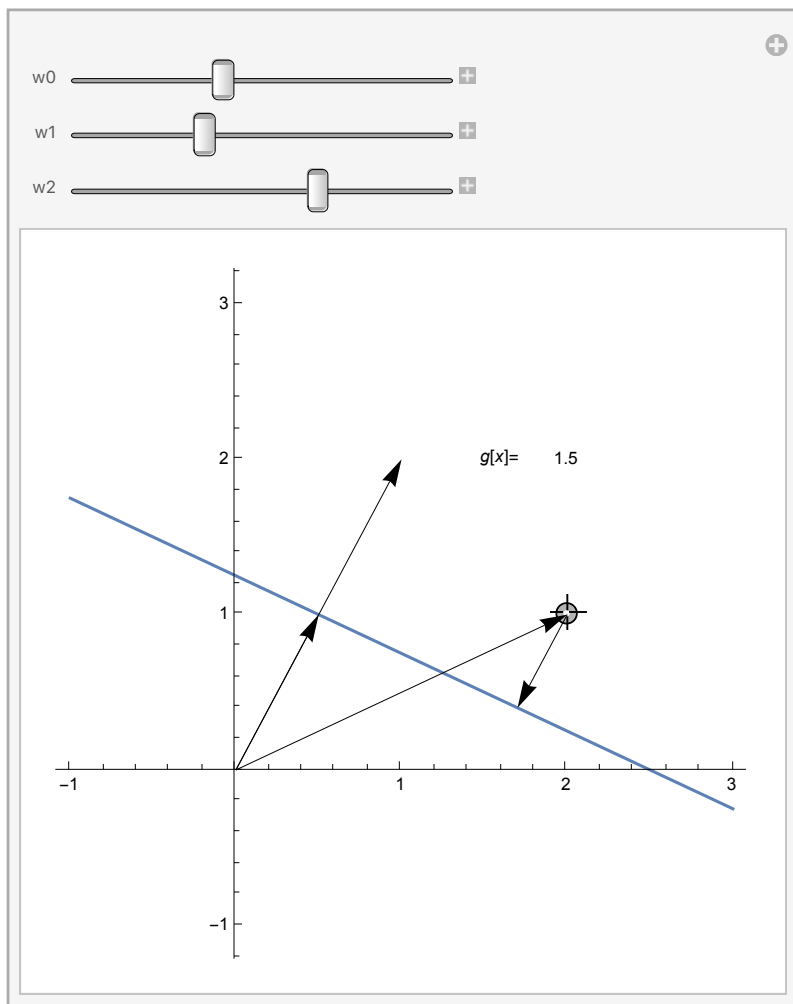


```

In[719]:= Clear[w1, w2, w, w0, x2, x1];
Manipulate[
  (* x0={2,1};*)
  w = {w1, w2}; wn = w / Norm[w];
  g[x_] := {w1, w2}.x + w0;
  gg = Plot[Tooltip[
    x2 /. Solve[{w1, w2}.{x1, x2} + w0 == 0, x2], "discriminant"], {x1, -1, 3}];
  ggg = Graphics[g[Dynamic[MousePosition["Graphics"]]]];
  Show[{ggg, Graphics[Inset["g[x]=", {1.6, 2}]],
    Graphics[Inset[ToString[g[x0]], {2, 2}]], Graphics[
      {Tooltip[Arrow[{{0, 0}, w}], "w"], Tooltip[Arrow[{{0, 0}, (-w0 / Norm[w]) * wn}],
        "-w0/|w|"], Tooltip[Arrow[{{0, 0}, x0}], "x"],
      Tooltip[Arrow[{x0, x0 - wn * g[x0] / Norm[w]}], "g(x)/|w|"}],
    PlotRange -> {{-1, 3}, {-1, 3}}, AxesOrigin -> {0, 0}, Axes -> True,
    AspectRatio -> 1, ImageSize -> Medium],
  {{w0, -2.5}, -6, 3}, {{w1, 1}, 0, 3}, {{w2, 2}, 0, 3},
  {{x0, {2, 1}}, Locator}]

```

Out[720]=

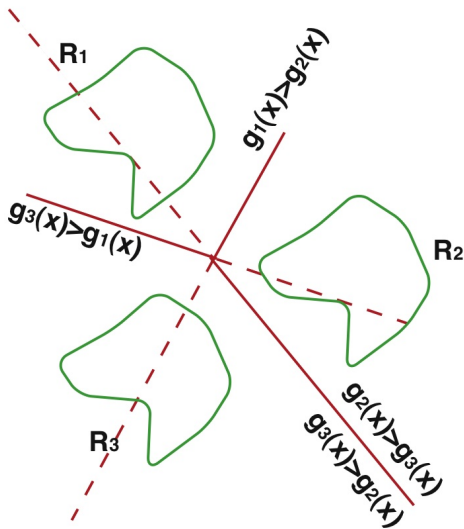


Multiple classes

Suppose there are c classes. There are a number of ways to define multiple class discriminant rules. One way that avoids undefined regions is:

$$g_i(\mathbf{x}) = \mathbf{w}_i \cdot \mathbf{x} + w_{i0}, \quad i = 1, \dots, c$$

Assign \mathbf{x} to the i th class if: $g_i(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq i$.



(Adapted from Duda & Hart, 1973)

It can be shown that this classifier partitions the input space into simply connected convex regions. This means that if you connect any two feature vectors belonging to the same class by a line, all points on the line are in the same class. Thus this linear classifier won't be able to handle problems for which there are disconnected clusters of features that all belong to the same class. Also, from a probabilistic perspective, if the underlying generative probability model for a given class has multiple peaks, this linear classifier won't do a good job either.

Fisher linear discriminant

Basic idea: find a projection that minimizes the within-class variance, while maximizing the between-class variance.

Task-dependent Dimensionality Reduction

Motivation

Later, when we consider the problem of dimensionality reduction, we will look at unsupervised methods, in particular principal components, PCA. But here the idea will be to find hyperplanes onto which we can project our input data, and from there divide up the hyperplane into decision regions. The idea is that the original input space may be impractically huge, but if we can find a subspace (hyperplane) that preserves the distinctions between categories as well as possible, we can make our decisions in smaller space. We will derive the Fisher linear "discriminant".

This is closely related to the psychology idea of finding "distinctive" features. E.g. consider bird identification. If I want to discriminate cardinals from other birds in my backyard, I can make use of the fact that (males) cardinals may be the only birds that are red. So even though the image of a bird can have lots of dimensions, if I project the image on to the "red" axis, I can do fairly well with just one number. How about male vs. female human faces?

A conceptually similar idea was used to learn diagnostic features in complex images (cf. Ullman et al., Hegde et al.).

We'll look at this more later when we learn about how to discover feature hierarchies.

Fisher's linear discriminant: building intuitions

Generative model: two nearby gaussian classes

Define two bivariate base distributions

```
In[721]:= (ar = {{1, 0.99}, {0.99, 1}};
ndista = MultinormalDistribution[{0, -1}, ar];)
(br = {{1, .9}, {.9, 2}};
ndistb = MultinormalDistribution[{0, 1}, br];)
```

Find the expression for the probability distribution function of ndista

```
pdf = PDF[ndista, {x1, x2}]
1.12822 e1/2 (-x1 (50.2513 x1 - 49.7487 (x2+1)) - (x2+1) (50.2513 (x2+1) - 49.7487 x1))
```

Use Mean[] and CovarianceMatrix[ndista] to verify the population mean and the covariance matrix of ndista

```
In[723]:= Mean[ndistb]
```

```
Out[723]= {0, 1}
```

```
In[724]:= Covariance[ndista]
```

```
Out[724]= {{1, 0.99}, {0.99, 1}}
```

- ▶ 3. Try different covariant matrices. Should they be symmetric? Constraints on the determinant of ar, br?
Make a contour plot of the PDF ndista

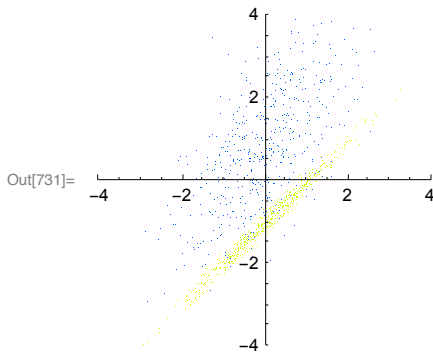
```
In[725]:= pdfa = PDF[ndista, {x1, x2}];
(*ContourPlot[pdfa, {x1, -3, 3}, {x2, -3, 3}, PlotPoints -> 64, PlotRange -> All] *)
```



```

In[726]:= nsamples = 500;
a = Table[RandomVariate[ndista], {nsamples}];
ga =
  ListPlot[a, PlotRange → {{-4, 4}, {-4, 4}}, AspectRatio → 1, PlotStyle → Hue[0.2`]];
b = Table[RandomVariate[ndistb], {nsamples}];
gb =
  ListPlot[b, PlotRange → {{-4, 6}, {-4, 4}}, AspectRatio → 1, PlotStyle → Hue[0.6`]];
Show[
  ga,
  gb]

```



Use `Mean[]` to find the *sample* mean of `b`. What is the sample *covariance* of `b`?

```
In[732]:= Mean[b]
```

```
Out[732]= {-0.0131852, 0.852518}
```

```
In[733]:= Covariance[b]
```

```
Out[733]= {{1.0251, 0.944284}, {0.944284, 2.12984}}
```

Try out different projections of the data by varying the slope (m) of the discriminant line

The basic idea is fairly straight forward, and if you are familiar with signal detection theory and the definition of d' , it should look familiar. We want to find a projection for which the squared difference between the means of the projections relative to the variance

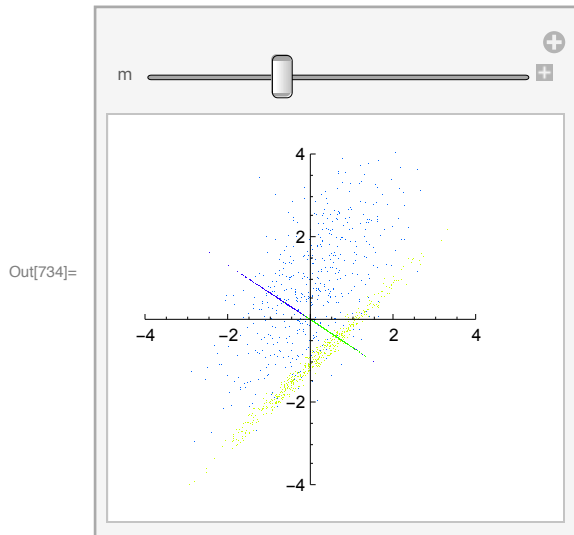
$$\frac{(\mu_b - \mu_1)^2}{\sigma_a^2 + \sigma_b^2}$$

is biggest.

```

In[734]:= Manipulate[
  wnvec = {1, m} / Sqrt[1 + m^2];
  aproj = (#1 wnvec &) /@ (a.wnvec);
  gaproj = ListPlot[aproj, AspectRatio -> 1, PlotStyle -> Hue[0.3]];
  bproj = (#1 wnvec &) /@ (b.wnvec);
  gbproj = ListPlot[bproj, AspectRatio -> 1, PlotStyle -> Hue[0.7]];
  Show[ga, gb, gbproj, gaproj, AspectRatio -> Automatic],
  {{m, -.66}, -2, 2}]

```



- 4. By trial and error, find a value of m that separates the classes well along the projection line. Plot out the marginal distributions relative to this line.

Calculate the "signal-to-noise" ratio along the projection line. Do this by taking the difference between the means divided by the square root of the product of the standard deviations along the line.

Theory and program demo for simple 2-class case

Here's the derivation of the Fisher discriminant for the 2-class case (see Duda and Hart for general case).

A measure of the separation between the projections is the difference between the means:

$$| w \cdot (m_a - m_b) |$$

and

$$m_a = \frac{1}{N} \sum_{i=1}^N x_i, \text{ summed over the } N \text{ } x \text{'s from class a}$$

$$m_b = \frac{1}{M} \sum_{i=1}^M x_i, \text{ summed over the } M \text{ } x \text{'s from class b}$$

where w (**wnvec**) is the unknown unit vector along the discriminant line .

In our case above, the vector difference between the means is:

```
In[735]:= Mean[a] - Mean[b]
```

```
Out[735]:= {0.0268344, -1.84623}
```

and the difference between the means projected onto a discriminant line is:

```
In[736]:= wnvec = {1, 2 / 3} / Sqrt[1 + 2 / 3 ^ 2];
```

```
In[737]:= wnvec . (Mean[a] - Mean[b])
```

```
Out[737]:= -1.08904
```

To improve separation, we can't just scale w , because the noise scales too.

We'd like the difference between the means to be large relative to the variation for each class. We can define a measure of the scatter for the projected samples in say class a ($a=1$), by:

$$\sum_{y \in \text{class } a} (y - \tilde{m}_a)^2$$

where \tilde{m}_a is the sample mean of the points from class a projected onto discriminant line and $y = \mathbf{w} \cdot \mathbf{x}$
Or in terms of the Mathematica example:

```
In[738]:= Apply[Plus, aproj - wnvec . Mean[a]];
```

The total scatter S is defined by the sum of the scatters for both classes (a and b).

$$S = \sum_{y \in \text{class } a} (y - \tilde{m}_a)^2 + \sum_{y \in \text{class } b} (y - \tilde{m}_b)^2$$

If we divide the above number by the total number of points, we have an estimate of the variance of the combined data along the projected axis.

We now have the basic ingredients behind the intuition for the Fisher linear discriminant. We'd like to find that \mathbf{w} for which J :

$$J(\mathbf{w}) = \frac{|\tilde{m}_a - \tilde{m}_b|^2}{S} = \frac{|\mathbf{w} \cdot (\mathbf{m}_a - \mathbf{m}_b)|^2}{S}$$

is biggest. We want to maximize the difference between the projected class means, while minimizing the dispersion of the data on the projected line.

One can show that $S = \mathbf{w}^T \cdot \mathbf{S}_W \cdot \mathbf{w}$, where

S_W is measure of within-class variation called the *within-class scatter matrix*:

$$S_W = \sum_{i=1}^2 \sum_{x \in \text{Class } i} (x - m_i) (x - m_i)^T$$

For the numerator, a measure of between class variation is the *between-class scatter matrix*:

$$S_B = (\mathbf{m}_1 - \mathbf{m}_2) \cdot (\mathbf{m}_1 - \mathbf{m}_2)^T$$

and the difference between the projected means can be show to be:

$$|\tilde{m}_a - \tilde{m}_b|^2 = \mathbf{w}^T \cdot S_B \cdot \mathbf{w}$$

Find \mathbf{w} (corresponding to slope) to maximize the criterion function

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \cdot S_B \cdot \mathbf{w}}{\mathbf{w}^T \cdot S_W \cdot \mathbf{w}}$$

The answer is given by:

$$w = S_W^{-1} \cdot (m_a - m_b)$$

Demo: Program to finding Fisher's linear discriminant

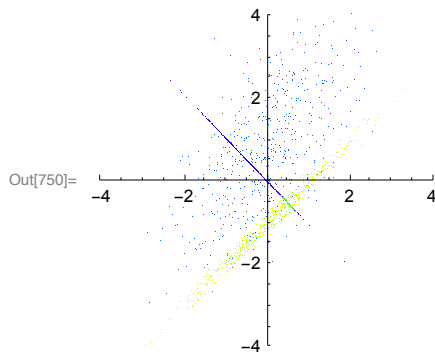
```
In[739]:= normalize[x_] := x / Sqrt[x.x];
```

```
In[740]:= ma = Mean[a];
mb = Mean[b];
```

```
In[742]:= Sa = Sum[Outer[Times, a[[i]] - ma, a[[i]] - ma], {i, 1, nsamples}];
Sb = Sum[Outer[Times, b[[i]] - mb, b[[i]] - mb], {i, 1, nsamples}];
Sw = Sa + Sb;
wldf = normalize[Inverse[Sw].(ma - mb)]
```

```
Out[745]= {0.69234, -0.721572}
```

```
In[746]:= aproj = (#1 wldf &) /@ (a.wldf);
gaproj = ListPlot[aproj, AspectRatio -> 1, PlotStyle -> Hue[0.3`]];
bproj = (#1 wldf &) /@ (b.wldf);
gbproj = ListPlot[bproj, AspectRatio -> 1, PlotStyle -> Hue[0.7`]];
Show[ga, gb, gbproj, gaproj]
```

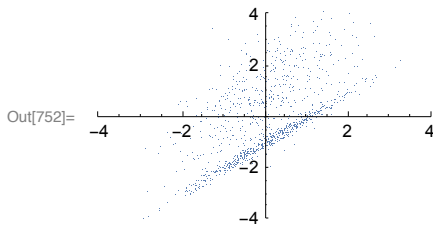


We started off with a 2-dimensional input problem and turned it into a 1-D problem. For the n-dimensional case, see Duda and Hart.

Compare with the principal component axes

Later we will look at principal components analysis as an unsupervised method to reduce dimensionality. But it is useful at this point to get an intuition to compare the two different ways to reduce dimensionality. The first component of PCA finds the dimension (axis) for which if you project the data on to it, the projections have the biggest variance.

```
In[751]:= c = Join[a, b];
ListPlot[c, PlotRange -> {{-4, 4}, {-4, 4}}]
```

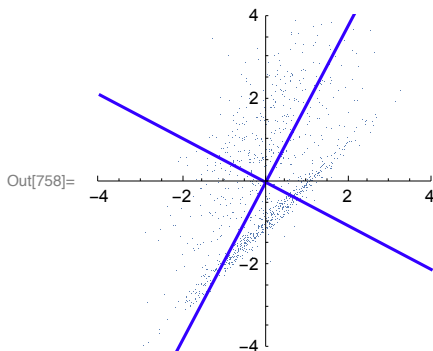


```
In[753]:= g1 = ListPlot[c, PlotRange -> {{-4, 4}, {-4, 4}}, AspectRatio -> 1];
```

```
In[754]:= auto = Covariance[c];
eigvalues = Eigenvalues[auto];
eigauto = Eigenvectors[auto];
```

```
In[757]:= gPCA = Plot[{eigauto[[1,2]]/eigauto[[1,1]] x,
eigauto[[2,2]]/eigauto[[2,1]] x},
{x,-4,4}, AspectRatio->1,
PlotStyle->{RGBColor[.2,0,1]}];
```

```
In[758]:= Show[g1, gPCA]
```



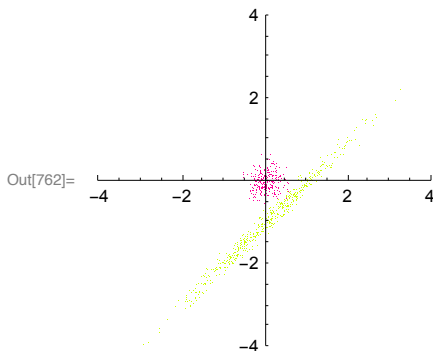
- ▶ 5. How does the principal component (biggest variance) compare with the Fisher discriminant line?
- ▶ 6. For classification is the Fisher discriminant always better than PCA?

If the main information for discrimination is in the variance, rather than the means, PCA can be better.

```

In[759]:= ndistc = MultinormalDistribution[{0, 0}, {{.05, 0}, {0, .05}}];
c = Table[RandomVariate[ndistc], {nsamples / 2}];
gc =
  ListPlot[c, PlotRange → {{-3, 3}, {-3, 3}}, AspectRatio → 1, PlotStyle → Hue[0.9]];
Show[
  ga,
  gc]

```



Support vector machines (SVMs) and kernel methods

Initialize

Read in SVM package:

Make sure the SVM package is downloaded in the default directory

```
In[763]:= SetDirectory[NotebookDirectory[]]
```

```
Out[763]= /Users/kersten/Sites/kersten-lab/courses/Psy5038WF2016/Lectures/Lect_18
          _SupervisedMLandNN
```

```
In[764]:= << MathSVMv7`
```

Sidenote: Some other useful directory functions:

```
In[765]:= FileNames[];
Directory[];
```

We assume a simple perceptron TLU to classify vector data \mathbf{x} into one of two classes depending on the sign of $g(\mathbf{x})$:

$$\text{decision}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b).$$

Given $g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$, recall that $g(\mathbf{x})/|\mathbf{w}|$ is the distance of a data point \mathbf{x} from a plane defined by $g(\mathbf{x}) = 0$. In support vector machines, the goal is to find the separating plane (i.e. find \mathbf{w} and b) that is as far as possible from any of the data points. The intuition is that this will minimize the probability of making wrong classifications when given new data at some future point.

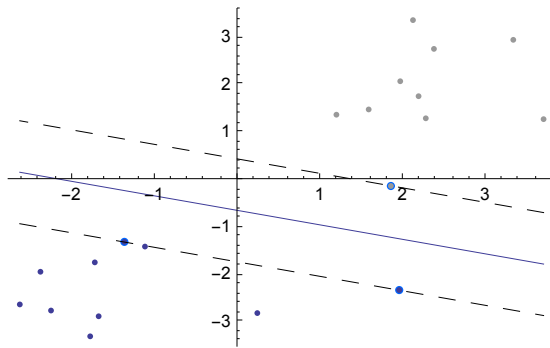
Let $\Pi_{\mathbf{w},b}$ represent the points on a potential decision line, and $d(\Pi_{\mathbf{w},b}, \mathbf{x}_i)$ is the distance between data

point x_i and the line. Formally, we want to solve:

$$\max_{(w,b)}(\min_i d(\Pi_{w,b}, x_i)), \quad (1)$$

where $d(\Pi_{w,b}, x_i) = g(x_i) / \|w\|$, i.e.

$$d(\Pi_{w,b}, x_i) = g(x_i) / \|w\| = |w \cdot x_i + b| / \|w\|$$



Two-class data (black and grey dots), their optimal separating hyperplane (continuous line), and support vectors (circled in blue). This is an example output of the `SVMPLOT` function in *MathSVM* (Nilsson et al., 2006). The width of the “corridor” defined by the two dotted lines connecting the support vectors is the called the *margin* of the optimal separating hyperplane. If the data are not linearly separable, the margin can be “softened” to allow some points to fall inside the corridor. The **soft-margin classifier** uses a regularizing term (similar to a prior on the parameters) to allow for a trade-off between classification error and tolerance.

The Primal Problem

It can be shown that the optimal separating hyperplane solving (1) can be found as the solution to the equivalent optimization problem

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

subject to $y_i(w^T x_i + b) \geq 1$,

Where y_i is the class label, either 1 or -1, for the i th point. Typically, equality will hold for a relatively small number of the data vectors. These data are termed *support vectors*. The solution (w, b) depends only on these specific points, and in effect contain all the information for the decision rule. The “dual problem”.

A simple linear SVM Example

Here's a demo of a simple SVM problem. It uses the add on package *MathSVMv7* written by Nilsson et al. (2006)

```
In[767]:= Clear[X, y];
len = 40;
X = Join[
  RandomReal[NormalDistribution[-2, 1], {len / 2, 2}],
  RandomReal[NormalDistribution[2, 1], {len / 2, 2}]];
y = Join[Table[1, {len / 2}], Table[-1, {len / 2}]];
```

We use the simple SVM formulation provided in *MathSVM* by the `SeparableSVM` function.

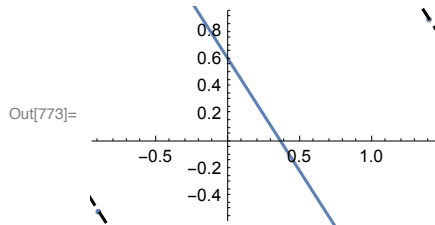
```
In[771]:=  $\tau = 0.01;$ 
```

```
 $\alpha = \text{SeparableSVM}[X, y, \tau]$ 
```

```
Out[772]:= {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0.276583, 0, 0,
0, 0, 0, 0, 0.276583, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0}
```

In the output figure below, the solid line marks the optimal hyperplane, and dotted lines mark the width of the corridor that joins support vectors (highlighted in blue).

```
In[773]:= SVMPlot[ $\alpha$ , X, y]
```



A Nonlinear Example: Using Kernels

What if the data are not linearly separable? The essential idea is to map the data (through some non-linear mapping, e.g. polynomial) to a higher-dimensional "feature" space to find the optimal hyperplane separating the data. The dot product gets replaced by a non-linear kernel function. For example, the polynomial kernel is given by:

```
In[774]:=  $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j)^d$ 
```

```
Out[774]=  $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \cdot \mathbf{x}_j)^d$ 
```

If $d = 1$, we have the standard dot product, but for $d = 2, 3$, etc.. we have polynomial functions of the elements of the vectors x . See Nilsson et al (2009), and paper by Jäkel (2009) for more information on kernels.

A demo for an application for nonlinear classification.

We'll use second-degree kernel, e.g.:

```
In[775]:= PolynomialKernel[{xx1, xx2, xx3}, {1, -1, -1}, 2]
```

```
Out[775]=  $(1 + xx1 - xx2 - xx3)^2$ 
```

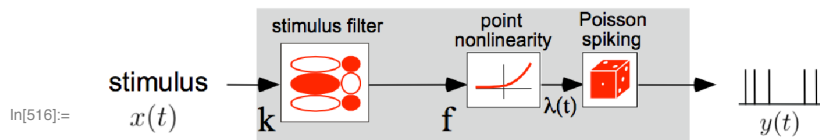
Some synthetic data which is not linearly separable.

Logistic regression, the generalized linear model, and the generic neuron model

This section is not about learning per se, but about how to motivate the class of functions one might like to fit.

The single (weight) layer perceptron is a special case of logistic regression, which in turn is a special case of a *generalized linear model*. Logistic regression has been studied by statisticians since the 1950s. And the *generalized linear model* is a basic tool in theoretical modeling of spiking neurons, in particular the linear-Nonlinear-Poisson model (LNP) is a special case, the first part of which should look familiar. In the figure below, the set of weights specifying the receptive field of a neuron are given by k :

Linear-Nonlinear-Poisson model



$$\text{conditional intensity (spike rate)} \quad \lambda(t) = f(k \cdot x(t))$$

(Figure due to Jonathan Pillow, lecture in Methods in Computational Neuroscience, 2011.)

Here's a general mathematical motivation for modeling the output as a probability of binary events.

The starting point is that we want to model a binary response Y (think action potential on or off) by the conditional probability:

$$P(Y=1 | X=x),$$

where X are the input features. Let's write

$$P(Y=1 | X=x) = p(x;w),$$

where $p(x;w)$ is a likelihood function of x for some parameters w ("k" in the above figure). How to model this function p ?

As we've done before, let's assume the simplest: try a linear function. But over what function of p ?

$p(x) = w \cdot x + b$ has problems because the left side is bounded between 0 and 1, and the right side is unbounded.

$\log p(x) = w \cdot x + b$ may seem nice because adding input features multiplies probabilities, but is bounded on only one side. A solution is to model log odds (of a spike vs. no spike) as a linear function:

$$\text{Log} \frac{p(x)}{1-p(x)} = w \cdot x + b$$

and after solving for $p(x)$,

$$p(x) = \frac{1}{1+e^{-(w \cdot x + b)}}$$

Looking familiar? This says that the probability of $Y=1$ is determined by taking a linear weighted sum of the inputs plus b (think bias) and deciding “yes” if this is bigger than zero. This is a linear classifier. This is also the probabilistic update rule we used for the Boltzmann machine.

To make things look even more familiar let’s calculate the expected value of $Y=1$ (of a neuron firing). First note that we can be a bit clever and write the probability of $Y= y_i$, whose values are 1 or 0 in a concise summary expression as:

$$p(y_i | x) = \frac{e^{y_i(w \cdot x + b)}}{1+e^{(w \cdot x + b)}}$$

Then the probability of firing, i.e. $y_i = 1$ is:

$$p(y_i = 1 | x) = \frac{1}{1+e^{-(w \cdot x + b)}} = \sigma(w \cdot x + b), \text{ where as we've seen before, } \sigma \text{ is sigmoidal.}$$

By the definition of expectation, the average rate is:

$$\sum_{y_i=0}^1 y_i p(y_i | x) = 0 \times p(y_i = 0 | x) + 1 \times p(y_i = 1 | x) = \sigma(w \cdot x + b)$$

This is stage 1 and 2 of our original generic neuron model.

The logistic regression model can be extended to **multiple class decisions**.

And as introduced above, these models in turn can be viewed as special cases of **generalized linear models**.

Note: The general linear model is different! It really is linear, and can be viewed as a special case. The generalized linear model is in general not. For applications of the generalized linear model to computational neuroscience, see Pillow (2007).

Naive Bayes and other methods

Naive Bayes. Assume we have a class variable C that takes on discrete values corresponding to labels. And assume a set of features $\{F_i\}$. Naive Bayes assumes that the probabilities of features are all conditionally independent of the class they came from:

$$p(C|F_1, \dots, F_n) = \frac{1}{Z} p(C) \prod_{i=1}^n p(F_i|C)$$

(You can view this as an extension to the cue integration problem studied earlier.) To do classification, we need a decision rule, such as MAP, which chooses the C that maximizes the posterior.

$$\operatorname{argmax}_C p(C | F_1, \dots, F_n)$$

http://en.wikipedia.org/wiki/Naive_Bayes_classifier

And others, AdaBoost, <http://en.wikipedia.org/wiki/AdaBoost>, and Random Forest, This Random Forest link also has a concise description of k-NN. http://en.wikipedia.org/wiki/Random_forest

Mathematica's Classify[] function demo

Mathematica tries to make it easy for non-experts to use its machine learning functions. This is convenient for learning concepts, and for many problems. The downside is that it can be difficult to find out exactly what the routines are doing.

Here we adapt code from one of *Mathematica*'s **Classify[]** examples and apply it to the noisy Xor problem. Define clusters sampled from normal distributions:

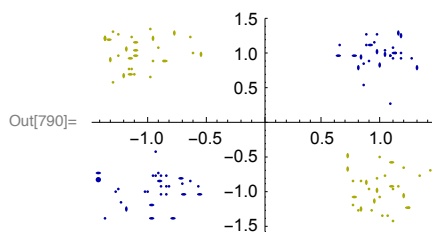
Set up training pairs

```
In[785]:= std = 0.05;
sampledata[center_] :=
  RandomVariate[MultinormalDistribution[center, std * IdentityMatrix[2]], 30];
```

Generate clusters, assign labels to each and plot;

```
In[787]:= clustercenters = {{-1, -1}, {-1, 1}, {1, -1}, {1, 1}}
clusters = sampledata /@ clustercenters;
colors = {Blue, Yellow, Yellow, Blue};
pl2 = ListPlot[clusters, PlotStyle -> Darker@colors]
```

```
Out[787]= {{-1, -1}, {-1, 1}, {1, -1}, {1, 1}}
```



```
In[791]:= Dimensions[clusters]
```

```
Out[791]= {4, 30, 2}
```

The finagling in the next line shows how to assign all of the 2D points to their corresponding colors

```
trainingdata0 = Flatten[Thread[Transpose[clusters][[#]] -> colors] & /@
  Range[Dimensions[clusters][[2]], 1];
```

```
In[793]:= Short[trainingdata0]
```

```
Out[793]/Short= {{-0.809456, -0.928514} -> ■, <<118>>, {1.27002, <<19>>} -> ■}
```

Set up the classifier, and classify

Now run the training pair associations through the Classifier, to automatically find the method, and do classifications on test data.

```
In[794]:= c30 = Classify[trainingdata0]
```

```
Out[794]= ClassifierFunction[ Method: NearestNeighbors  
Number of classes: 2]
```

What method did it pick? With the default, setting Mathematica decides.

How to calculate probabilities and make decisions:

```
In[795]:= c3[{-2, 3}, "Probabilities"]
```

```
Out[795]= < |  → 0.435751,  → 0.564249 | >
```

► 7. Get information on method, parameters, etc.:

```
ClassifierInformation[c3, "MethodDescription"];
ClassifierInformation[c3, #] & /@ {"NeighborsNumber", "DistanceFunction"};
ClassifierInformation[c3];
ClassifierInformation[c3, "Properties"];
```

► 8. Change method and parameters:

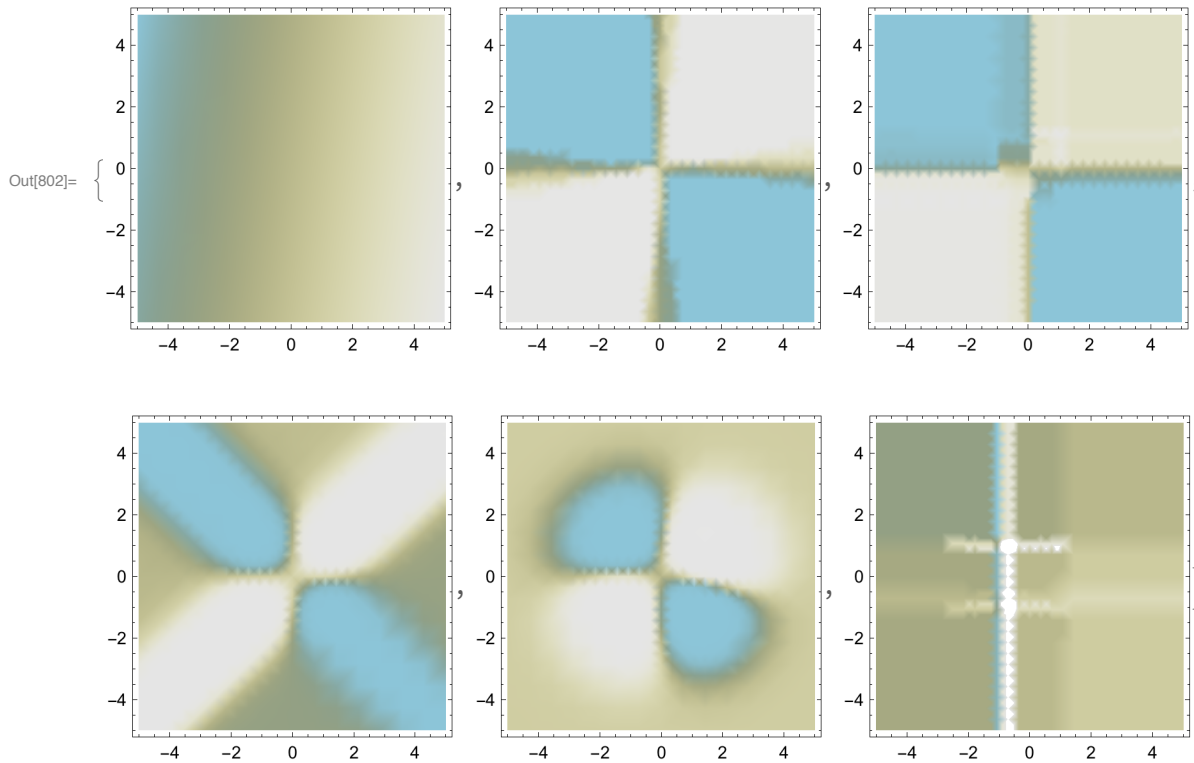
```
In[800]:= ctemp = Classify[trainingdata0, Method → {"NearestNeighbors", "NeighborsNumber" → 4}]
```

```
Out[800]= ClassifierFunction[ Method: NearestNeighbors  
Number of classes: 2]
```

Visualize decision probabilities

```
In[801]:= plotprob[method_] := Module[{},
  c3 =
    Classify[Flatten[Thread[Transpose[clusters][[#]] → colors] & /@ Range[30], 1],
      Method → method];
  DensityPlot[Normal@c3[{x, y}, "Probabilities"][[1]], {x, -5, 5},
    {y, -5, 5}, Exclusions → None, ColorFunction → "LightTerrain"]
];
```

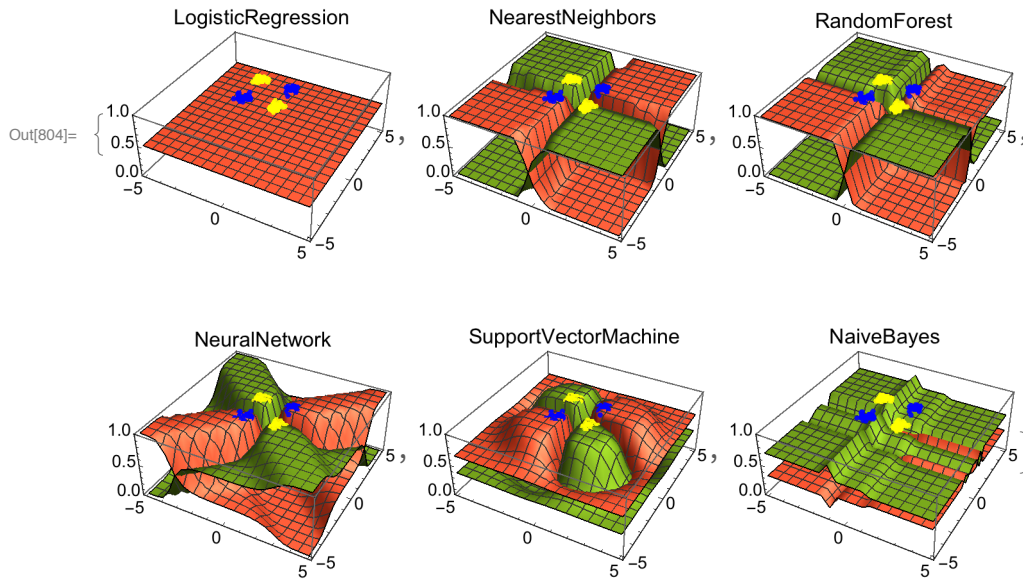
```
In[802]:= plotprob /@ {"LogisticRegression", "NearestNeighbors",
  "RandomForest", "NeuralNetwork", "SupportVectorMachine", "NaiveBayes"}
```



```
In[803]:= plotprobabilities[method_] := Module[{c, proba},
  c3 = Classify[Flatten[
    Thread[Transpose[clusters][[#]] → colors] & /@Range[30], 1], Method → method];
  proba = Table[Append[{x, y}, #] & /@Lookup[c3[{x, y}, "Probabilities"], colors],
    {x, -5, 5, .5}, {y, -5, 5, .5}];
  proba = Flatten[#, 1] & /@Transpose[proba, {3, 2, 1, 4}];
  Show[
    ListPlot3D[proba, PlotRange → {0, 1}, PlotLabel → method],
    ListPointPlot3D[Map[Append[#, 1] &, clusters, {2}],
      PlotStyle → ({Opacity[.9], #} & /@ colors)], ImageSize → 150
  ]
];
```

Visualize decision probabilities using 3D plots

```
In[804]:= plotprobabilities /@ {"LogisticRegression", "NearestNeighbors",
  "RandomForest", "NeuralNetwork", "SupportVectorMachine", "NaiveBayes"}
```



- ▶ 9. How do the various classification methods compare with respect to interpolation vs. extrapolation?
- ▶ 10. Try a noisier XOR by setting `std = 0.25`. This increases the probability of overlapping samples.

Next time

Doing inference by sampling

Transitioning to Python

Why?

Pros:

Free

Rapid growth in packages for scientific computation, including spiking neural networks (<http://briansimulator.org>), machine learning, image processing, computer vision, and psychophysics (PsychoPy).

Lots of packages, examples, tutorials. More specialized packages than *Mathematica*. For example, deep convolutional networks: <http://deeplearning.net/tutorial/lenet.html>.

Cons:

Lots and lots of packages, examples, tutorials. Packages are a set of moving targets that are may live or fade away and die. Not as many specialized packages as matlab. Installation can be a pain.

Scientific python: `numpy`, `scipy`, ...

Raw python is too basic. We need add on modules. So in order of importance, we will always start by loading `numpy` for creating and manipulating numerical arrays, `scipy` for scientific functions (e.g. opti-

