# Lecture 9: Supervised learning

STAT598z: Intro. to computing for statistics

Vinayak Rao

Department of Statistics, Purdue University

## Supervised learning

We are given training data $(X, Y) = \{(x_1, y_1), \cdots, (x_N, y_N)\}$

- X: independent variables, inputs, predictors, features
- Y: dependent variables, outputs, response

$x \in R^P$ (usually)

- regression: $y \in R$
- classification: $y \in \{0, 1\}$
- structured prediction: More complicated high-dimensional spaces with dependent components (e.g. the space of images or sentences)

We assume $y_i = f(x_i) + \varepsilon_i$

$\varepsilon$ is noise (includes randomness and approximations)

- Independently and identically distributed (i.i.d.) according to some probability distribution (e.g. the Gaussian)

Given the training set $(X, Y)$, we want to estimate $f$:

- to study the relation between x and y
- to make predictions of y's for unobserved x's

Good predictors can be hard to interpret

## Parametric learning

Index functions $f$ by a finite-dimensional parameter vector E.g. linear regression

- Parameters are coefficients of a hyperplane
- Parameters have a clear interpretation
- Can be a bad approximation of reality

# Linear regression

via the `lm` function in R

```
In [ ]:  DataIncm <- read.table('Data/Income2.csv',header=T,sep=',')
         fit <- lm(Income ~ Education, DataIncm)
```

The first argument is a formula

- takes the form response ~ predictors
- response is a linear combination of predictors
- above we have just one predictor: $Education$
- $Income = \beta_1 \cdot Education + \beta_0 + \epsilon$

Second argument unnecessary if variables in formula exist in current environment

See documentation for other optional arguments

Can print `fit`:

```
fit
```

This is not all the information in `fit` (why?)

- Try `typeof()`, `class()`, `str()`
- Try plotting it

Observe fit contains the entire dataset!

Can disable with `model = FALSE` option

Can directly plot with ggplot :

In [ ]:
```r
library('ggplot2')
plt1 <- ggplot(DataIncm, aes(x=Education, y = Income)) +
        geom_point(size=2, color='blue') +
        theme(text=element_text(size=10))
```

In [ ]:
```r
plt1 + geom_smooth(method='lm', se=FALSE, #Disable std. errors
                   color='magenta', size=2)
```

Can regress against Seniority

In [ ]: 
```
fit <- lm(Income ~ Seniority, DataIncm)
```

Can regress against both Education and Seniority

In [ ]: 
```
fit <- lm(Income ~ Education + Seniority, DataIncm)
```

- + does *not* mean input is sum of Educ. and Sen.

Rather: $Income = \beta_2 \cdot Seniority + \beta_1 \cdot Education + \beta_0 + \varepsilon$

For the former, use I:

fit <- lm(Income ~ I(Education + Seniority), DataIncm)

- $Income = \beta_1 \cdot (Seniority + Education) + \beta_0 + \varepsilon$

## Prediction

```
In [ ]:  fit <- lm(Income ~ Education + Seniority, DataIncm)
         fit
```

How do we make predictions at a new set of locations? E.g. (15, 60) and (20, 160)?

```
In [ ]:  pred_locn <- data.frame(Education=c(15,20), Seniority= c(60,160))
         predict.lm(fit, pred_locn)
```

```
In [ ]: edu_pred <- 10:25
        sen_pred <- seq(0,200,10)
        pred <- data.frame(Education=rep(edu_pred, length(sen_pred)),
                     Seniority=rep(sen_pred, each=length(edu_pred) ))
        p_val <- predict.lm(fit, pred)
        pred_full <- cbind(pred,p_val)
```

```
In [ ]: plt <- ggplot(DataIncm, aes(x=Education, y=Seniority,
                            color=Income))+
           geom_tile(data=pred_full, aes(x=Education, y=Seniority,
                                    color=p_val, fill=p_val)) +
           geom_point(size=1) + theme(text=element_text(size=10)) +
           scale_color_continuous(low='blue', high='orange') +
           scale_fill_continuous(low='blue', high='orange') +
           geom_point(shape=1,size=1,color='black') +
             guides(fill=FALSE)
```

```python
In [ ]: plt
```

Specifying a model for `lm`

| Symbol | Meaning | Example |
|--------|---------|---------|
| + | Include variable | x + y |
| : | Interaction between vars | x + y + z + x:z + y:z |
| * | Variables and interactions | (x + y) * z |
| ^ | Vars and intrcns to some order | (x + y + z)^3 |
| - | Delete variable | (x + y + z)^3 - x:y:z |
| poly | Polynomial terms | poly(x,3) + (x + y) * z |
| I | New combination of vars | I(x*y + z) |
| 1 | Intercept | x - 1 |

See documentation and http://ww2.coastal.edu/kingw/statistics
/R-tutorials/formulae.html (http://ww2.coastal.edu/kingw/statistics
/R-tutorials/formulae.html)

# Generalized linear model

A linear model with Gaussian noise is often inappropriate. E.g.

- response is always positive
- count valued response
- {0, 1} or binary-valued as in classification

A better model might be:

$$response = g(\sum_{i=1}^{N} \beta_i \cdot predictor_i) + \varepsilon$$

$g$ is a 'link' function, $\varepsilon$ is no longer Gaussian

Can fit in R with `glm()` (see documentation)