# Practical Advice on Machine Learning Projects

## Not included in (most) textbooks

MF1633055

杨润琦

# General Process

- Problem Setup
- Data Exploration
- Data cleaning
- Feature Construction
- Feature Transformation
- Feature selection
- Architecture Design
- Model Tuning
- Error Analysis

" Any intelligent fool can make things bigger and more complex. It takes a touch of genius, and a lot of courage, to move in the opposite direction. "
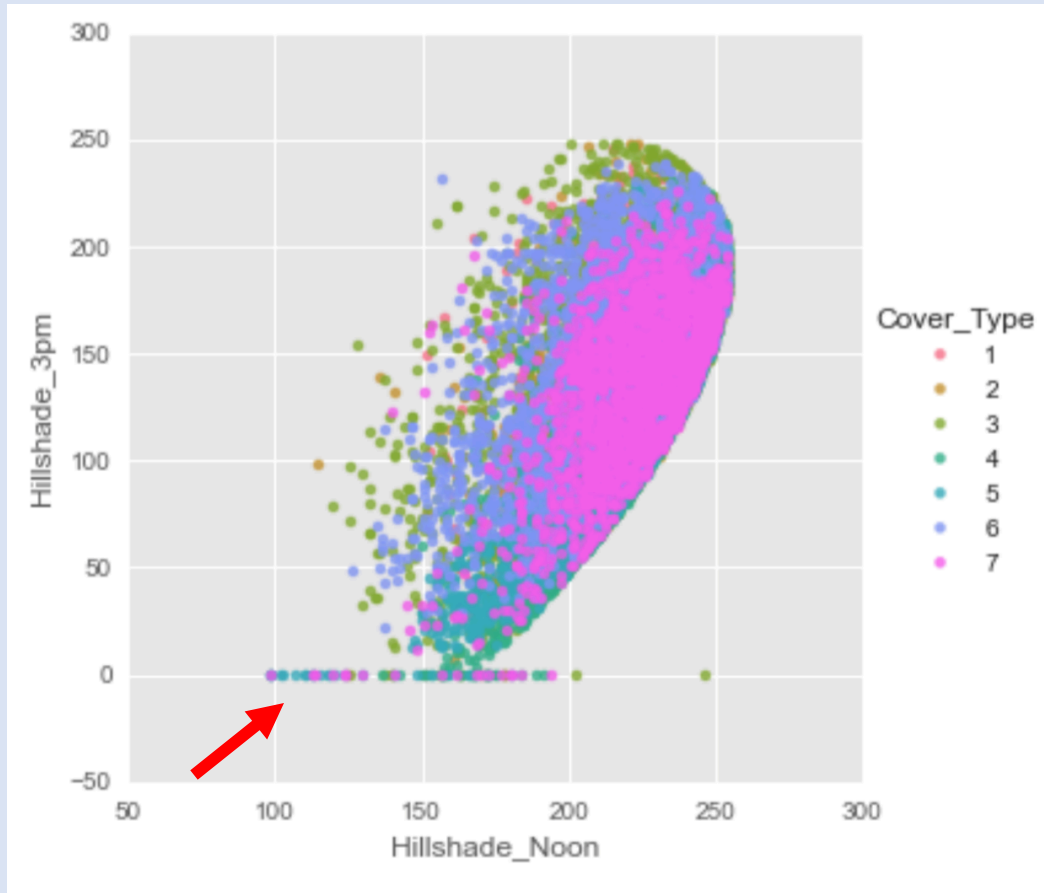
Before the competition begins….

- Identify the problem
- Collect the appropriate data
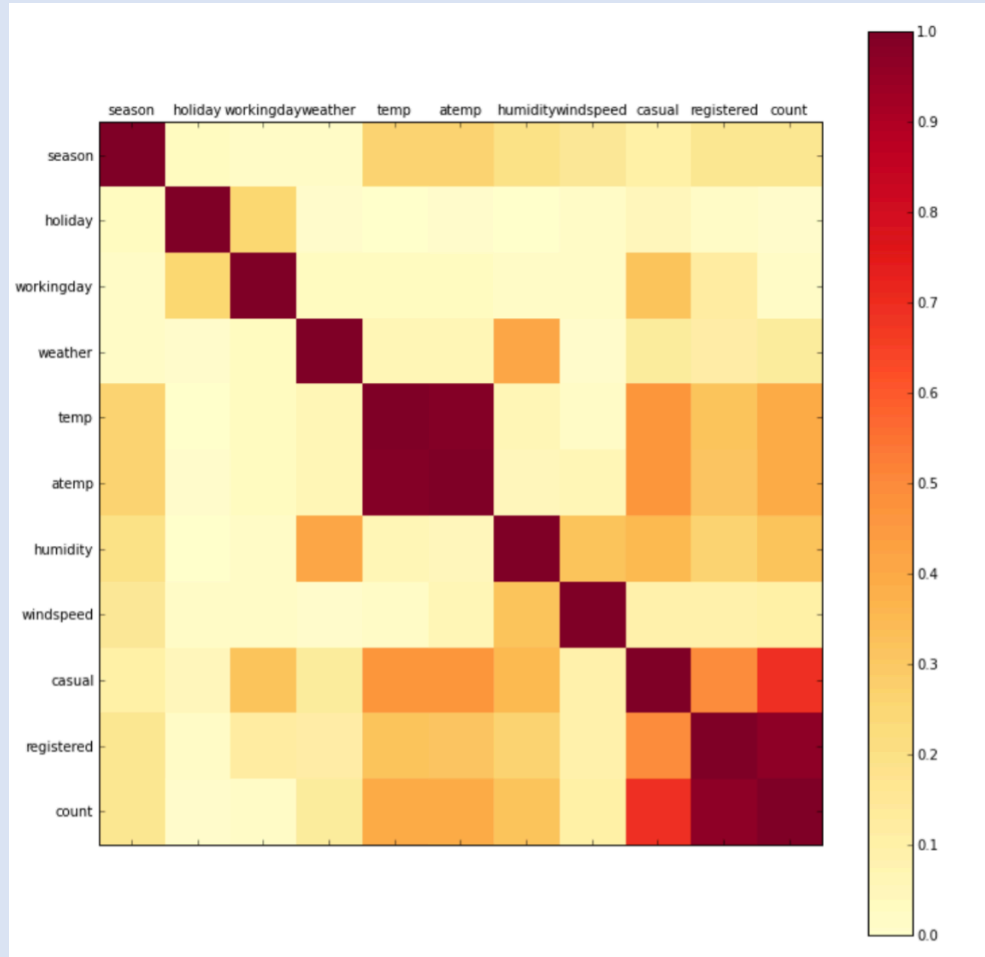- split the data into training and test datasets

# No standard way, but...



```
sns.jointplot(x='attr1', y='attr2', data=X,
kind='scatter', hue='label')
```

- Problem Setup
- **Data Exploration**
- Data cleaning
- Feature Construction
- Data Transformation
- Feature selection
- Architecture Design
- Model Tuning
- Error Analysis

# No standard way, but…



```
sns.heatmap(data.corr())
```

- Problem Setup
- **Data Exploration**
- Data cleaning
- Feature Construction
- Data Transformation
- Feature selection
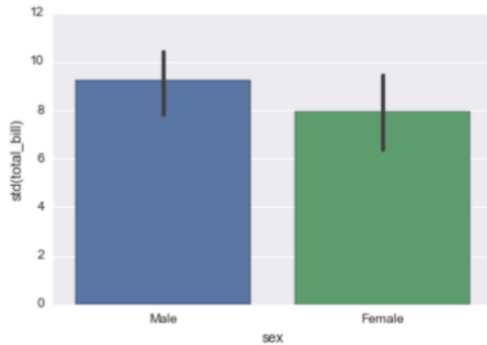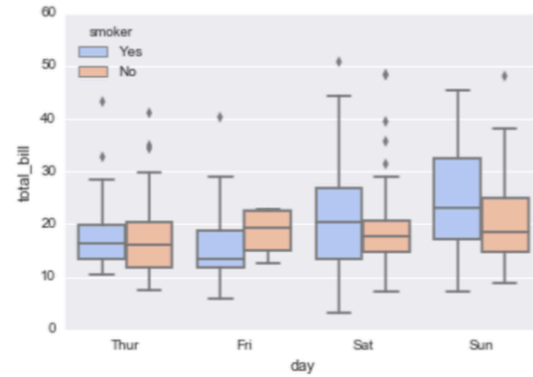- Architecture Design
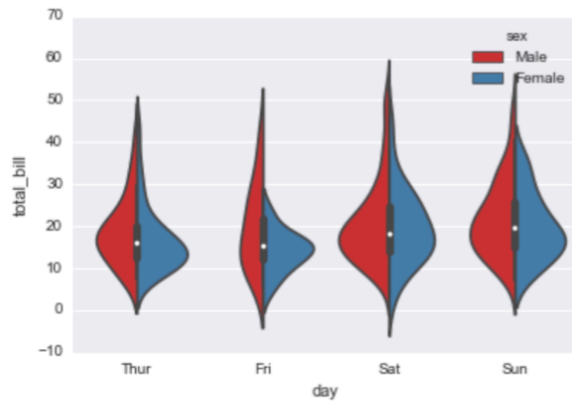- Model Tuning
- Error Analysis

# No standard way, but...
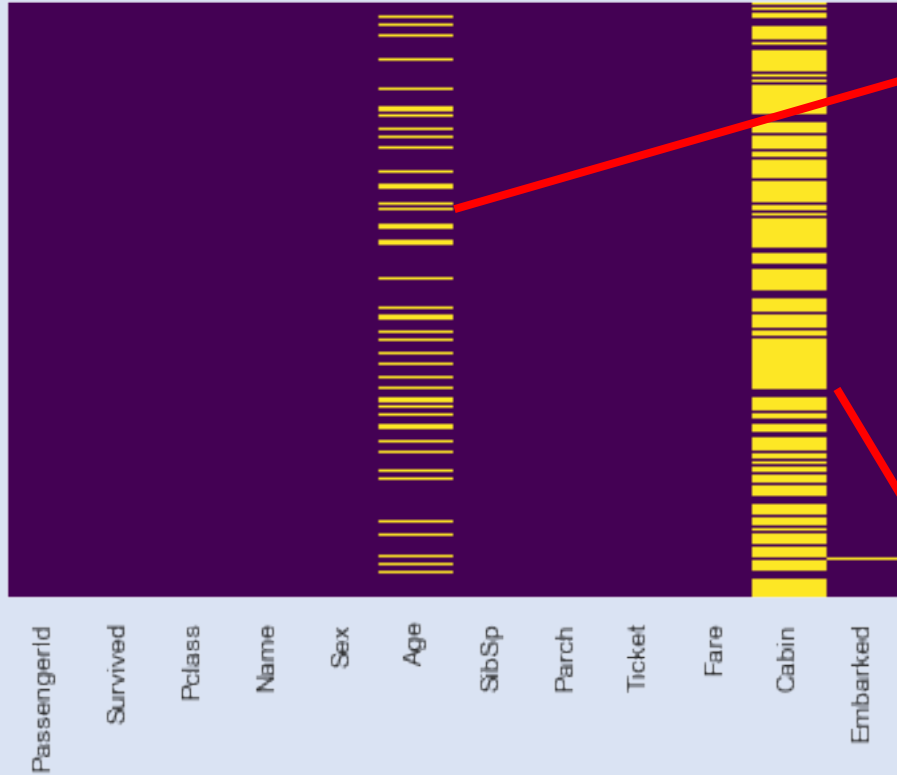


barplot (countplot)

boxplot

violinplot

stripplot (swarmplot)

- Problem Setup
- **Data Exploration**
- Data cleaning
- Feature Construction
- Data Transformation
- Feature selection
- Architecture Design
- Model Tuning
- Error Analysis

# Fill in missing values



`sns.heatmap(X.isnull())`

- Use attribute mean, tiered average…
- Build a model
  - test/dev data can be included

- Drop the column
- change to another feature: known/unknown

Once you spot a trend/pattern in data exploration, try to convert it to a feature.

Normalization

Discretization

Binarization

- Problem Setup
- Data Exploration
- Data cleaning
- Feature Construction
- **Data Transformation**
- Feature selection
- Architecture Design
- Model Tuning
- Error Analysis

- Problem Setup
- Data Exploration
- Data cleaning
- Feature Construction
- Data Transformation
- **Feature selection**
- Architecture Design
- Model Tuning
- Error Analysis

Generated using data from Sogou User Grouping Contest, 2016

Ensemble
- voting
- stacking

Clustering
- train a model in each cluster
- use cluster label as a new feature

Much more…

# Grid Search

```
1  from sklearn.grid_search import GridSearchCV
2  param_grid = {
3      'C': [0.1, 1, 10, 100, 1000],
4      'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
5      'kernel': ['rbf']}
6  grid = GridSearchCV(SVC(), param_grid, refit=True, verbose=3)
7  grid.fit(X_train,y_train)
```
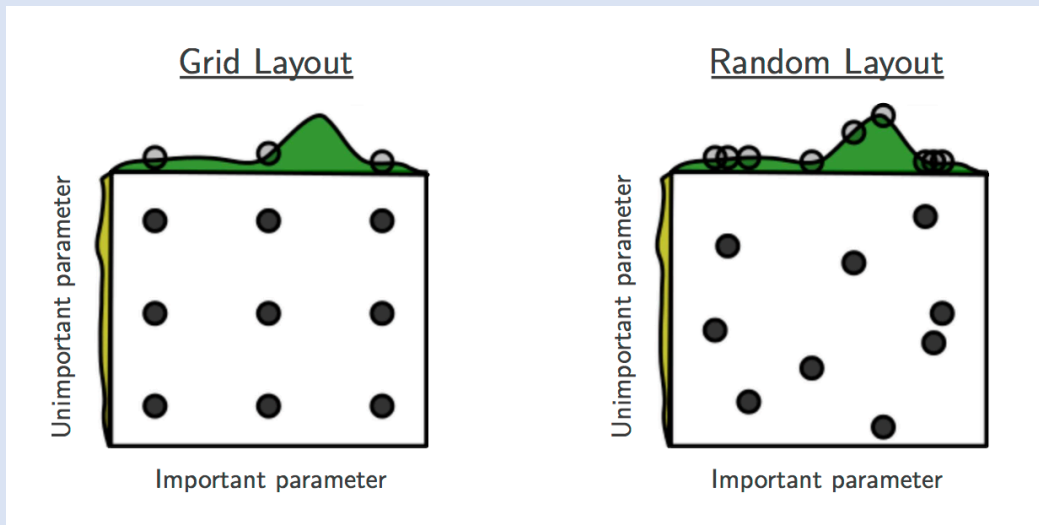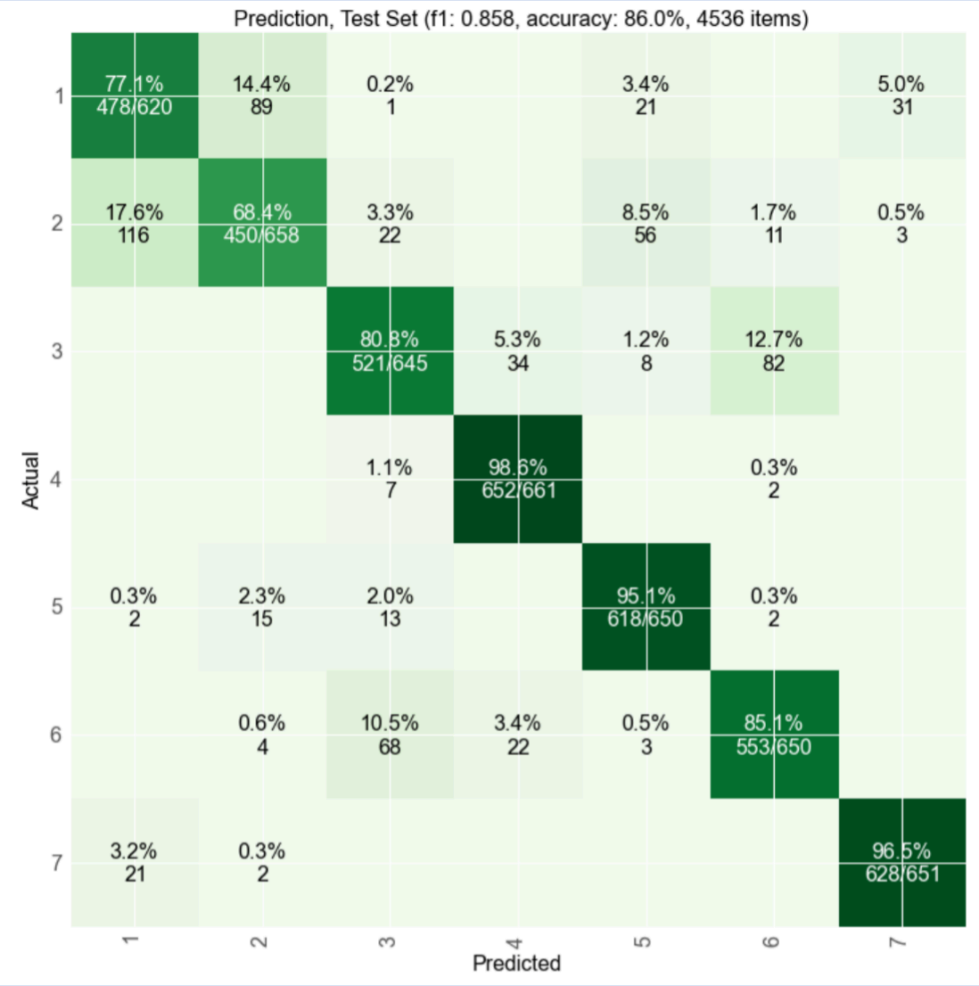
# Random Search



- Problem Setup
- Data Exploration
- Data cleaning
- Feature Construction
- Data Transformation
- Feature selection
- Architecture Design
- **Model Tuning**
- Error Analysis

# Classification
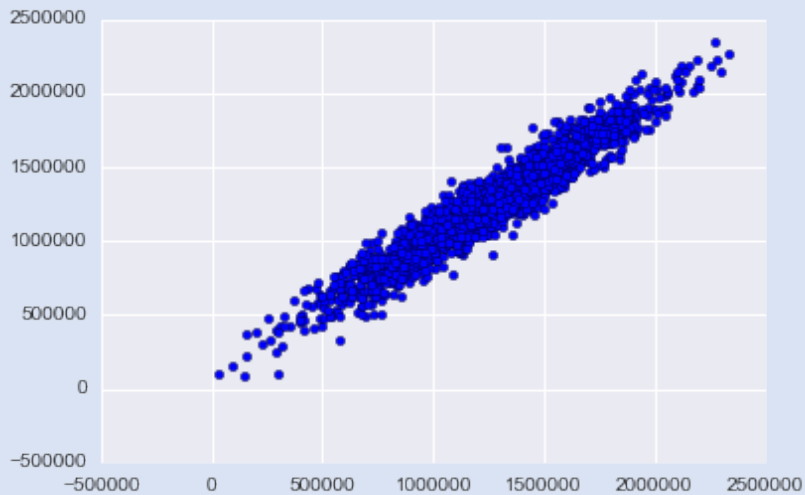# Matrix Plot of Confusion Matrix



- Problem Setup
- Data Exploration
- Data cleaning
- Feature Construction
- Data Transformation
- Feature selection
- Architecture Design
- Model Tuning
- **Error Analysis**
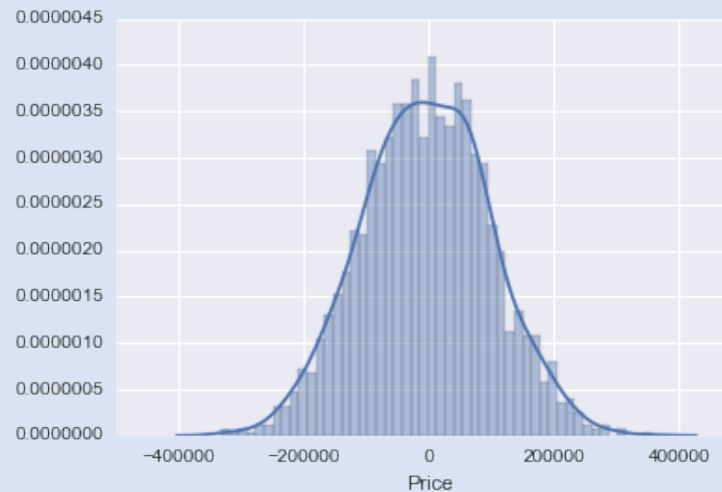
# Matrix Plot of Confusion Matrix (Code)

```python
cm = confusion_matrix(y_true, y_pred, labels=labels)
cm_sum = np.sum(cm, axis=1, keepdims=True)
cm_perc = cm / cm_sum * 100
annot = np.empty_like(cm).astype(str)
nrows, ncols = cm.shape
for i in range(nrows):
    for j in range(ncols):
        c = cm[i, j]
        p = cm_perc[i, j]
        if i == j:
            s = cm_sum[i]
            annot[i, j] = '%.1f%%\n%d/%d' % (p, c, s)
        elif c == 0:
            annot[i, j] = ''
        else:
            annot[i, j] = '%.1f%%\n%d' % (p, c)
cm = pd.DataFrame(cm, index=labels, columns=labels)
cm.index.name = 'Actual'
cm.columns.name = 'Predicted'
sns.heatmap(cm, annot=annot, fmt='')
```

# Regression



```
plt.scatter(y_dev, predictions)
```



```
sns.distplot(
    (y_dev-predictions),
    bins=50)
```

- Problem Setup
- Data Exploration
- Data cleaning
- Feature Construction
- Data Transformation
- Feature selection
- Architecture Design
- Model Tuning
- **Error Analysis**

Q&A

Thanks for your time!

MF1633055

杨润琦