

دليل المفاهيم الرياضية والإحصائية المستخدمة في مشروع التنبؤ بمضاعفات السكري

المواضيع الرياضية والإحصائية المستخدمة بالترتيب

1. الإحصاء الوصفي (Descriptive Statistics)

المفاهيم المستخدمة:

- المقاييس المركزية:
 - Mean (المتوسط الحسابي)
 - Median (الوسيط)
 - Mode (النوال)
- مقاييس التشتت:
 - Standard Deviation (الانحراف المعياري)
 - Variance (التباين)
 - Range (المدى)
- الإحصائيات الأساسية:
 - Min/Max (أصغر وأكبر قيمة)
 - Quartiles (الأرباع)
 - Percentiles (المئينات)

التطبيق في المشروع

```
df.describe() # الإحصائيات الوصفية الأساسية
df['Age'].mean() # متوسط العمر
df['HbA1c'].std() # الانحراف المعياري لمستوى السكر
```

2. تصور البيانات (Data Visualization)

المفاهيم المستخدمة:

- Histograms:** عرض توزيع المتغيرات المستمرة
- Box Plots:** مقارنة التوزيعات واكتشاف القيم الشاذة
- Scatter Plots:** دراسة العلاقات بين المتغيرات
- Bar Charts:** عرض المتغيرات الفئوية

التطبيق في المشروع

```
plt.hist(df['Age']) # توزيع الأعمار
sns.boxplot(x='Gender', y='HbA1c') # مقارنة مستوى السكر حسب الجنس
```

3. معالجة البيانات المفقودة (Missing Data Handling)

المفاهيم المستخدمة

- **Mean Imputation:** تعويض القيم المفقودة بالمتوسط
- **Mode Imputation:** تعويض القيم المفقودة بالمنوال
- **Pattern Analysis:** تحليل نمط البيانات المفقودة

التطبيق في المشروع

```
df['Age'].fillna(df['Age'].mean()) # تعويض بالمتوسط
df['Gender'].fillna(df['Gender'].mode()[0]) # تعويض بالمنوال
```

4. ترميز المتغيرات (Variable Encoding)

المفاهيم المستخدمة

- **Label Encoding:** تحويل المتغيرات النصية إلى أرقام
- **Feature Scaling:** توحيد مقاييس المتغيرات
- **StandardScaler:** التطبيع باستخدام المتوسط والانحراف المعياري

التطبيق في المشروع

```
from sklearn.preprocessing import LabelEncoder, StandardScaler
le = LabelEncoder()
scaler = StandardScaler()
```

5. اختبار الفرضيات وتحليل العلاقات

المفاهيم المستخدمة

- **Cross-tabulation:** الجداول المتقاطعة
- **Group Comparisons:** مقارنة المتوسطات بين المجموعات
- **Distribution Analysis:** تحليل التوزيعات

التطبيق في المشروع

```
pd.crosstab(df['Gender'], df['Complications']) # الجدول المتقاطع
df.groupby('Complications')['Age'].mean() # متوسط العمر حسب المضاعفات
```

6. تقسيم البيانات (Data Splitting)

المفاهيم المستخدمة

- **Train-Test Split:** (تقسيم البيانات 80% تدريب، 20% اختبار)
- **Stratified Sampling:** الحفاظ على توزيع المتغير التابع
- **Random State:** ضمان إعادة الإنتاج

التطبيق في المشروع:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
stratify=y)
```

7. خوارزميات التعلم الآلي (Machine Learning Algorithms)

(أ) **(Tree-based Models)** النماذج القائمة على الأشجار

Random Forest:

- **Ensemble Learning مفهوم:** دمج عدة نماذج
- **Bootstrap Aggregating (Bagging):** أخذ عينات عشوائية
- **Feature Importance:** قياس أهمية المتغيرات
- **Out-of-Bag Error:** قياس الخطأ

Gradient Boosting:

- **Boosting Algorithm:** تعلم تتابعي
- **Gradient Descent:** تحسين الدالة
- **Loss Function:** دالة الخسارة

(ب) **(Linear Models)** النماذج الخطية

Logistic Regression:

- **Sigmoid Function:** $\sigma(z) = 1/(1 + e^{(-z)})$
- **Maximum Likelihood Estimation:** تقدير الاحتمالية العظمى
- **Log-Odds:** اللوغاريتم الطبيعي للاحتمالات

(ج) **(Instance-based)** نماذج المسافة

Support Vector Machine (SVM):

- **Kernel Methods:** دوال النواة
- **Margin Optimization:** تحسين الهامش
- **Support Vectors:** النقاط الداعمة

8. التحقق المتقاطع (Cross-Validation)

المفاهيم المستخدمة:

- **K-fold Cross-Validation:** أجزاء K تقسيم البيانات لـ
- **Stratified K-fold:** الحفاظ على توزيع الهدف
- **Mean and Standard Deviation:** متوسط وانحراف معياري النتائج

التطبيق في المشروع:

```
from sklearn.model_selection import cross_val_score
cv_scores = cross_val_score(model, X_train, y_train, cv=5)
```

9. مقاييس التقييم (Evaluation Metrics)

المفاهيم الرياضية:

Accuracy (الدقة):

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

Precision (الدقة الموجبة):

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

Recall/Sensitivity (الحساسية):

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

F1-Score:

$$\text{F1} = 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$$

حيث:

- TP = True Positives (الإيجابيات الصحيحة)
- TN = True Negatives (السلبات الصحيحة)
- FP = False Positives (الإيجابيات الخاطئة)
- FN = False Negatives (السلبات الخاطئة)

10. منحنيات الأداء (Performance Curves)

المفاهيم المستخدمة:

ROC Curve (منحنى خاصية التشغيل المتلقي):

- **True Positive Rate (TPR):** معدل الإيجابيات الصحيحة
- **False Positive Rate (FPR):** معدل الإيجابيات الخاطئة
- **AUC (Area Under Curve):** المساحة تحت المنحنى

Precision-Recall Curve:

- مفيد للبيانات غير المتوازنة
- يركز على الأداء في الفئة الإيجابية

11. أهمية المتغيرات (Feature Importance)**المفاهيم المستخدمة:**

- **Tree-based Feature Importance:** قياس إسهام كل متغير
- **Statistical Ranking:** ترتيب المتغيرات إحصائياً
- **Gini Impurity:** مقياس عدم النقاء في الأشجار

📖 مصادر التعلم المقترحة**(المستوى الأساسي (البداية****1. كتب أساسية:**

- **"The Elements of Statistical Learning"** - Hastie, Tibshirani, Friedman
 - يغطي: الإحصاء، التعلم الآلي، النظرية الرياضية
- **"Introduction to Statistical Learning with R"** - James, Witten, Hastie, Tibshirani
 - يغطي: مقدمة عملية للتعلم الآلي مع التطبيق

2. دورات أونلاين مجانية:

- **Khan Academy: Statistics and Probability**
 - الرابط: <https://www.khanacademy.org/math/statistics-probability>
 - يغطي: الإحصاء الوصفي، الاحتمالات، اختبار الفرضيات
- **Coursera: "Machine Learning" by Andrew Ng**
 - يغطي: أساسيات التعلم الآلي، الخوارزميات الأساسية
- **edX: "Introduction to Probability and Statistics"**
 - يغطي: الاحتمالات، التوزيعات، الاستنتاج الإحصائي

(المستوى المتوسط (التطبيق**3. Python والتعلم الآلي:**

- **"Python for Data Analysis"** - Wes McKinney
 - يغطي: pandas, numpy, matplotlib, seaborn
- **"Hands-On Machine Learning"** - Aurélien Géron
 - يغطي: scikit-learn, tensorflow, التطبيق العملي

4. دورات متخصصة:

- **Coursera: "Applied Data Science with Python Specialization"**
 - التصور، التعلم الآلي، تحليل النصوص، الشبكات الاجتماعية، Python: دورات تغطي 5
- **DataCamp: Statistics and Machine Learning courses**
 - دورات تفاعلية عملية

(المستوى المتقدم (التعمق

5. مراجع متقدمة:

- **"Pattern Recognition and Machine Learning"** - Christopher Bishop
 - يغطي: النظرية الرياضية العميقة، الخوارزميات المتقدمة
- **"Machine Learning: A Probabilistic Perspective"** - Kevin Murphy
 - يغطي: النهج الاحتمالي، البايزي، النماذج المتقدمة

6. مواقع ومراجع متخصصة:

- **Scikit-learn Documentation**
 - الرابط: <https://scikit-learn.org/stable/>
 - مرجع شامل لجميع الخوارزميات
- **Towards Data Science (Medium)**
 - مقالات متخصصة في علم البيانات
- **StatQuest YouTube Channel**
 - شرح مبسط للمفاهيم الإحصائية والتعلم الآلي

ترتيب الدراسة المقترح

(المرحلة الأولى (4-6 أسابيع

1. ابدأ بالإحصاء الوصفي

- المصدر: Khan Academy Statistics
- التركيز: المتوسط، الوسيط، الانحراف المعياري، التوزيعات

2. للبيانات Python تعلم

- المصدر: Python for Data Analysis
- التركيز: pandas, numpy, matplotlib

(المرحلة الثانية (4-6 أسابيع

3. مقاييس التقييم والتحقق

- المصدر: Scikit-learn Documentation
- التركيز: Accuracy, Precision, Recall, F1-Score, Cross-validation

4. تصور البيانات

- المصدر: DataCamp Visualization courses

- فهم البيانات بصرياً, seaborn, plotly, التركيز

(المرحلة الثالثة (6-8 أسابيع

5. خوارزميات التعلم الآلي الأساسية

- المصدر: Introduction to Statistical Learning
- التركيز: Linear Regression, Logistic Regression, Decision Trees

6. خوارزميات متقدمة

- المصدر: Hands-On Machine Learning
- التركيز: Random Forest, SVM, Gradient Boosting

(المرحلة الرابعة (4-6 أسابيع

7. التطبيق العملي المتقدم

- المصدر: Coursera Applied Data Science
- التركيز: مشاريع حقيقية، معالجة البيانات، النشر

8. التقييم والتحسين

- المصدر: Elements of Statistical Learning
- التركيز: فهم عميق للنظرية، تحسين الأداء

💡 نصائح للدراسة الفعالة

1. النهج العملي:

- اطبق كل مفهوم فور تعلمه
- حقيقية للتدريب datasets استخدم
- اكتب الكود بنفسك بدلاً من النسخ

2. التدرج في التعلم:

- لا تتجاوز مفهوم قبل فهمه جيداً
- ارجع للأساسيات عند الحاجة
- اربط المفاهيم الجديدة بما تعلمته سابقاً

3. الممارسة المستمرة:

- حل مسائل وتمارين يومياً
- Kaggle شارك في مسابقات
- اعمل على مشاريع شخصية

4. المراجعة الدورية:

- راجع المفاهيم الأساسية شهرياً
- اعمل ملخصات لكل موضوع
- اختبر نفسك بانتظام

بعد كل مرحلة، يجب أن تكون قادراً على

المرحلة الأولى:

- Python حساب الإحصائيات الوصفية يدوياً وبـ
- فهم معنى كل مقياس إحصائي
- تحليل البيانات بصرياً

المرحلة الثانية:

- احترافية visualizations إنشاء
- تفسير مقاييس التقييم
- cross-validation تطبيق

المرحلة الثالثة:

- تطبيق خوارزميات مختلفة
- مقارنة أداء النماذج
- فهم متى تستخدم كل خوارزمية

المرحلة الرابعة:

- بناء نظام تنبؤ كامل
- تحسين أداء النماذج
- نشر النموذج للاستخدام

🔍 مراجع إضافية متخصصة

للرياضيات:

- **Linear Algebra:** "Linear Algebra and Its Applications" by Gilbert Strang
- **Calculus:** Khan Academy Calculus courses
- **Probability Theory:** "A First Course in Probability" by Sheldon Ross

للبرمجة:

- **Python:** "Automate the Boring Stuff with Python"
- **Pandas:** "Python for Data Analysis" official documentation
- **NumPy:** "NumPy User Guide"

للتطبيقات الطبية:

- **Medical Statistics:** "Medical Statistics at a Glance" by Aviva Petrie
- **Epidemiology:** "Modern Epidemiology" by Rothman, Greenland, Lash

□□□□□ □□□□□□ □□ □□□□□ □□□□□□ □□□□□ □□□□□□ □□□□ □□□□□□□□.

تاريخ الإنشاء: أغسطس 2025
آخر تحديث: أغسطس 2025